

AD-A100 652

GEORGIA INST OF TECH ATLANTA SCHOOL OF INFORMATION A-ETC F/8 9/2  
OPERATIONAL SURVIVABILITY IN GRACEFULLY DEGRADING DISTRIBUTED P-ETC(U)  
DEC 80 E W MARTIN N00014-79-C-0231  
61T-ICS-01/04 NL

UNCLASSIFIED

1 of 1  
AD-A100 652



END  
DATE  
FILMED  
7-81  
DTIC

AD A100652

LEVEL *12*

AD \_\_\_\_\_

TECHNICAL REPORT  
GIT-ICS-81/04

*118-049434*

## OPERATIONAL SURVIVABILITY IN GRACEFULLY DEGRADING DISTRIBUTED PROCESSING SYSTEMS

By  
Edith Waisbrot Martin

Prepared for  
OFFICE OF NAVAL RESEARCH  
800 N. QUINCY STREET  
ARLINGTON, VIRGINIA 22217

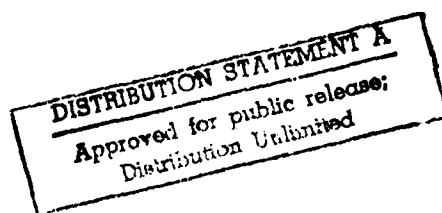


Under  
Contract No. N00014-79-C-0231  
GIT Project No. G36-636

December, 1980

**GEORGIA INSTITUTE OF TECHNOLOGY**  
SCHOOL OF INFORMATION AND COMPUTER SCIENCE  
ATLANTA, GEORGIA 30332

DTIC FILE COPY



1980



81 5 26 115

THE RESEARCH PROGRAM IN  
FULLY DISTRIBUTED PROCESSING SYSTEMS

OPERATIONAL SURVIVABILITY  
IN  
GRACEFULLY DEGRADING  
DISTRIBUTED PROCESSING SYSTEMS

TECHNICAL REPORT  
GIT-ICS-81/04

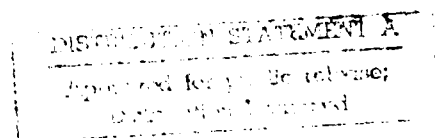
Edith Waisbrot Martin

December, 1980

Office of Naval Research  
800 N. Quincy Street  
Arlington, Virginia 22217

Contract No. N00014-79-C-0231  
GIT Project No. G36-636

The Georgia Tech Research Program in  
Fully Distributed Processing Systems  
School of Information and Computer Science  
Georgia Institute of Technology  
Atlanta, Georgia 30332



THE FINDINGS IN THIS REPORT ARE NOT TO BE CONSTRUED  
AS AN OFFICIAL DEPARTMENT OF THE NAVY POSITION,  
UNLESS SO DESIGNATED BY OTHER AUTHORIZED DOCUMENTS.

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

14 REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM	
1. REPORT NUMBER GIT-ICS-81/04	2. GOVT ACCESSION NO. AD-A100652	3. RECIPIENT'S CATALOG NUMBER	
4. TITLE (and Subtitle) Operational Survivability in Gracefully Degrading Distributed Processing Systems.		5. TYPE OF REPORT & PERIOD COVERED Technical Report	
7. AUTHOR(s) 10/ Waiskrist Dr. Edith M. Martin		6. PERFORMING ORG. REPORT NUMBER GIT-ICS-81/04	
9. PERFORMING ORGANIZATION NAME AND ADDRESS School of Information and Computer Science Georgia Institute of Technology Atlanta, Georgia 30332		8. CONTRACT OR GRANT NUMBER(s) N00014-79-0231	
11. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research 800 N. Quincy Street Arlington, Virginia 22217		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS C-	
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) 12/86		12. REPORT DATE 11 December 1980	
		13. NUMBER OF PAGES 155 + xi	
		15. SECURITY CLASS. (of this report) Unclassified	
		15a. DECLASSIFICATION DOWNGRADING SCHEDULE N/A	
16. DISTRIBUTION STATEMENT (of this Report) The United States Government is authorized to reproduce and distribute reprints for government purposes			
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report) DISTRIBUTION STATEMENT A Approved for public release; Distribution Unlimited			
18. SUPPLEMENTARY NOTES The findings in this report are not to be construed as an official Department of the Navy position unless so designated by other authorized documents.			
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) Distributed Data Processing      Multiple Regression Survivability Performance Evaluation Data Analysis			
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) Under this contract a simulator was designed and developed which will model possible distributed system network topologies, distributed system application topologies and their effect on application system performance as the con- figuration of the distributed system network is continuously and arbitrarily reduced. The objective of the model is to aid in development of a measure of survivability which can subsequently be used to evaluate and compare alternative distributed system designs for specific battlefield applications.			

DD FORM 1 JAN 73 1473

EDITION OF 1 NOV 55 IS OBSOLETE

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

410044

## SUMMARY

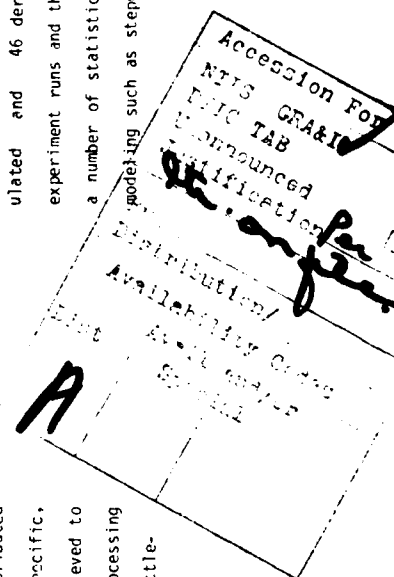
To date the concept of survivability as it pertains to distributed processing systems has been an intuitive one. The objective of this research is to present this concept quantitatively. Toward this end a number of hypotheses are presented, namely, that survivability must be measured in a nontrivial or indirect manner; that survivability is a function of a number of attributes, all of which are necessary to adequately explain or predict survivability; that the attributes which describe survivability are large in number and complex in interaction; and that because of these characteristics traditional performance, survivability and reliability measures are inadequate. This research proposes to demonstrate the applicability of standard experimental design and regression analysis techniques to the field of computer science in general and modeling of distributed systems in specific.

To test these hypotheses a computer model which supports the simulation of distributed processing systems for the purpose of evaluation and experimentation was constructed. This model, called SURSIM, models distributed system networks, application systems and several distribution/redistribution approaches and the effect of these on application system performance as the configuration of the distributed system network is continuously and arbitrarily reduced. In specific, the simulator represents and manipulates those attributes believed to be important in evaluating the survivability of distributed processing systems which must operate in real-time environments such as battle-

field situations. The controlled factors include; distributed system topology, network size (i.e., number of nodes); node processing, memory and communications capacity; applications system size, connectivity and interaction requirements; distribution strategies and extent of distributed system degradation.

SURSIM comprises modeling of the application system, distributed processing system, assignment of the application system to the distributed processing system, and mutation of the distributed system with subsequent reconfiguration of the application system. The capability to analyze application system performance based on application system requirements and distributed system capabilities is provided. In addition, SURSIM has the ability to implement degradation procedures which reduce software application system requirements to accommodate degraded distributed system capabilities. Considerable input and operational data logging is performed as the simulator is exercised. Tables and matrices describing the distributed and application systems being examined are output. Experimental results are generated in tabular and machine readable form to facilitate manual and computerized analysis.

A 2<sup>k</sup>-P Fractional Factorial experiment was conducted using the simulator as an experimental tool. One hundred and twenty-eight experimental cases were run in which 11 different factors were manipulated and 46 derived measures monitored. The results of the 128 experiment runs and the subsequent 300,000 subcases were examined using a number of statistical techniques. Several approaches to regression modeling such as stepwise and all possible subset regression were used



to build explanatory models from the results collected. Ten of these models proved to serve well in an explanatory capacity, consequently, data splitting was employed to assess the value of these models when used as predictors. It was determined that three models which function well in an explanatory role also serve well in predicting survivability and performance.

Thirty-two candidate regressors are used in identifying the 10 best subset models. The coefficients of these regressors are approximately equivalent in sign and magnitude across models. All variables remain proportional with the introduction and removal of other variables, thereby demonstrating extreme stability. The explanatory adequacy of models built using these variables is in all instances in excess of .8 which is very acceptable for a factor screening experiment. The adequacy in prediction of these models ranges between -.39 and +.11 with some models predicting very well and other predicting very poorly. By constructing satisfactory explanatory and predictive models, this research demonstrates that the concept of operational survivability and performance as proposed can be expressed quantitatively. Further, it is shown that major factors include the distributed system network, application system and distribution policy as initially hypothesized. Nine factors are found in all models. These are number of nodes in the distributed system, distributed system connectivity, module memory requirements, module to module interaction frequency, distribution policy, percent nodes lost, initial assignment results, available processing capacity at the end of the subcase and the interaction of all application related variables.

The research conducted here identifies the variables important to operational survivability and to some extent tells how large changes in these important variables affect the response. Future experimentation which provides either a large number of factor levels or finer granularity in possible variable values should permit greater resolution in the simulator results and their subsequent application. The results presented in this dissertation demonstrate the applicability of traditional experimentation and regression analysis in the field of computer science as well as the feasibility of measurements which can serve as measurements for distributed systems. The models developed represent a promising initial step in the quantification of operational survivability as it applies to gracefully degrading distributed processing systems.

## TABLE OF CONTENTS

SUMMARY . . . . .	Page v
LIST OF TABLES . . . . .	x
LIST OF ILLUSTRATIONS . . . . .	xi
CHAPTER	
I. INTRODUCTION . . . . .	1
II. BACKGROUND . . . . .	6
III. APPROACH . . . . .	12
IV. PROCEDURE . . . . .	20
V. SIMULATOR RESULTS . . . . .	39
VI. ANALYSIS PART I - PROCESS. . . . .	54
VII. ANALYSIS PART II - INTERPRETATION. . . . .	81
VIII. CONCLUSIONS AND RECOMMENDATIONS. . . . .	102
APPENDIX A DESCRIPTION OF DATA USED IN DESIGNED EXPERIMENTS . . . . .	115
APPENDIX B TEN OPTIMAL SUBSET MODELS . . . . .	126
APPENDIX C TEN MULTIPLE LINEAR REGRESSION MODELS BUILT FROM ESTIMATION SET DATA B . . . . .	139
ACKNOWLEDGEMENTS . . . . .	152
VITA . . . . .	153
BIBLIOGRAPHY. . . . .	154

## LIST OF TABLES

Table	Page
1. Experimental Factors and Factor Levels . . . . .	24
2. Experiment Factors and Pseudo-factors . . . . .	25
3. $2^{K-P}$ Fractional Factorial Design . . . . .	27
4. Interpretation of Example Treatment Combination . . . . .	31
5. Survivability Simulator Experiment Run Description . . . . .	46
6. Application System Topology Interaction Incidence Matrix . . . . .	47
7. Application System Requirements . . . . .	47
8. Degradation Policy . . . . .	48
9. Distributed System Topology Interaction Incidence Matrix . . . . .	49
10. Distributed System Capability . . . . .	49
11. Resources Remaining After Initial Assignment . . . . .	50
12. Module to Node Assignment . . . . .	50
13. Sample Data Log . . . . .	51
14. Statistics from BMDP Multiple Linear Regression Analysis . . . . .	65
15. Comparison of Models Constructed by Three Regression Methods . . . . .	71
16. Candidate Regressors . . . . .	75
17. X Factors Present in Ten Best Subset Models . . . . .	77
18. Comparison of Ten Models Fitted on Estimation Data Set . . . . .	78
19. Coefficients for Variables in Ten Best Subset Models . . . . .	83
20. Rank Ordering of 10 Best Subset Models . . . . .	89



# LIST OF ILLUSTRATIONS

Figure	Page
1. Index Tables for Four Level Factors . . . . .	30
2. Survivability Simulator Flow Diagram . . . . .	38
3. Multiple Linear Regression Model with All Candidate Repressor Variables . . . . .	64
4. Stepwise Regression Model . . . . .	67
5. Optimum Model According to Mallows Cp Criterion . . . . .	69
6. Optimum Model According to Adjusted $R^2$ Criterion . . . . .	70
7. Average Performance for Different Distributed System Topologies and Distribution Policies . . . . .	94

## CHAPTER I

### INTRODUCTION

#### Overview

The effectiveness of any distributed system design must be viewed against a backdrop of predetermined weights and priorities. To many users, the main benefit to be derived from the distributed approach to application system processing is increased capability to satisfy application system requirements despite the loss of a portion of the distributed system resources. The extent of that capability is herein termed "survivability." Inherent in the concept of survivability is that of "graceful degradation." Gracefully degrading systems are those which attempt to provide a high quality of service by reconfiguring the system or network or by reallocating resources when a fault is detected. This term is used to imply that performance may decrease with successive failures but it may not be catastrophically effected.

This research investigates the concept of survivability in gracefully degrading systems. It examines distributed system resources, processing nodes and associated links, which can be lost before a given application system required to execute on that distributed system must function in a degraded mode or experience failure. Whereas the determination of survivability has thus far been primarily judgemental based on a spectrum of performance variables, it is the intent of this research to express this concept quantitatively. Toward that end the applicability of standard experimental design methods and

regression analysis techniques to issues in performance evaluation are demonstrated.

Evaluation of survivability is of importance, for example, to the U.S. Army. Modern warfare has made automation on the battlefield essential to provide the commander with timely information on which to base his decisions and for weapon system and equipment control. Automation is required to enhance the speed, accuracy and dependability of battlefield systems that perform the functions of command, control, communications, intelligence, air defense, weapons control, surveillance, electronic warfare, sensor control, field artillery, navigation, logistics, and administration. Foremost among the goals for battlefield automation are operational effectiveness and continuity of operations. The ability to meet these goals is determined by 1.) the quality of the automation hardware and software components, 2.) the compatibility or interchangeability of these components, and 3.) the capability of the components to cooperate together to accomplish assigned tasks.

As an example, the Military Computer Family (MCF) program proposes to address each of these issues. The MCF program addresses quality of hardware and software components by utilizing the state-of-the-art instruction set architecture (ISA) and high order language (HOL) which appear to best fit the projected needs of Army automation systems. These were selected with extreme emphasis on potential reliability, performance and maintainability attributes. The MCF hardware and software members were chosen for their ability to meet the widest possible spectrum of Army automation system requirements and

thereby provide the foundation or standard for such systems. Via standardization the issue of interchangeability is accommodated. This research addresses the third segment of the operational effectiveness/continuity of operations duo, that is the capability of the MCF members to function together as a distributed system serving the requirements of specific Army application systems in battlefield situations.

#### Distributed System Design Considerations

In distributed systems the concern is for systems composed of many processing and memory components working together to serve a common application. For the most part designs are desired which provide decentralized control of the system, that is the controller does not reside in a single processor. Distributing application system tasks over a number of processing components can result in greater computing speed and capacity than is possible with a single processor of the same approximate cost. This benefit is achieved by customizing the selection and configuration of components to best match application system requirements, and thereby minimizing inefficient use of computer system resources. Distributed systems are believed to be inherently less costly to modify or upgrade because single, relatively small, components can be added or replaced rather than whole systems. Decentralization of resources and application system processing can yield additional benefits with respect to reliability and fault tolerance in that distribution of resources and activities can be so arranged such that the likelihood of single-point failures is reduced. There are numerous ways in which utilization of distributed systems can be advantageous. Currently, many advantages are obtained through

distribution over processing components which are physically close, i.e., within a one mile radius of one another. Distribution at more geographically separated locations entails additional complexity due to transmission delays and increased susceptibility of the distributed system to failures associated with communication losses or noise. Many applications, however, such as military real-time systems, require interaction with geographically dispersed system components. This dispersion may be for reasons inherent in the application or for purposes of system survivability or reduced vulnerability to loss of a portion of the distributed processing system locations. In this research the proximity of the processing components will be important as well as the qualification that the processing components be operating together to serve a single application.

Designing distributed processing systems to-date is not a well understood practice. This is true in part because the distributed system concept is still somewhat new and because it introduces additional variables and complexities into the design process. For example, distributed systems contain concurrent processes which must share resources and data without the benefit of centralized control, data and application system processes are distributed and possibly replicated at multiple locations, and communication management and protocols are often cumbersome and complex. Each of these design issues in the area of distributed processing systems appears to be more complex than its centralized counterpart. Thus far, no general design approach exists which adequately addresses the extremely complex and varied goals of distributed systems. Among the unresolved design

issues are database distribution and management, distributed control, task distribution, fault tolerance, and performance prediction.

## CHAPTER II

### BACKGROUND

There has been considerable research over the past twenty years in the area of fault tolerant computing. Initially the focus of that research was on the hardware of single processor systems. Fault tolerant computer investigations centered on models of ultra reliable systems having long mission time requirements such as those used in space exploration (19). These were typically uniprocessor systems with requirements for extremely large mean-time-to-first-failure (MTTF). Loss of a processing node was tantamount to catastrophic failure of the function or system served. The software support component of such systems was small and uncomplicated usually consisting of some minimum executive support required to execute the application software.

Later applications with large continuous processing demands presented a need for computer systems with high availability (3.4). The increased throughput and reliability required to support these applications was achieved by the introduction of a special class of multiple processor systems. These were repairable computer systems which embodied the concepts of redundancy and standby sparing (16.4). The important criterion of processor reliability shifted from mean-time-to-first-failure to mean-time-between-failures (MTBF). Support software was more complex than before. Software design had to address issues such as placement of software tasks and monitoring of hardware system components to detect failures and institute recovery

procedures. No matter what hardware backup scheme was employed, the reliability of the system became much more noticeably dependent upon the operation of the executive software. Software control for the most part was either centralized often realizing an underlying master slave relationship or functionally dispersed. As a consequence of this relationship, systems were still extremely vulnerable to single point failure.

Vulnerability to single point failure in part motivated development of distributed hardware and software systems which incorporated distributed control. Prior reliability and fault tolerance concepts laid the foundation for a new system reliability approach which attempts to provide a high quality of service by reallocating resources, i.e. reassigning tasks, or by reconfiguring the system or network, i.e. changing physical interconnection or routing algorithm, when a fault is detected. Such systems are termed "gracefully degrading" systems. This term is used to imply that performance may decrease with successive failures but it may not be catastrophically effected. Techniques for graceful degradation are particularly useful when applied to loosely coupled processor systems such as networks or fully distributed systems, that is, systems the components of which have high potential for autonomous operation. Many Command, Control and Communication, and tactical systems, for example, have significant impact if they experience catastrophic failure. The consequence of failure justifies the additional expense of hardware and software needed to allow military systems to withstand partial system failures.

Graceful degradation techniques are not in wide use currently. In part this is true because there is still a great lack of knowledge of the software organization required to implement graceful degradation. Also, there is a lack of adequate analytical models with which to evaluate such systems. Some research (2,17,22,12) in the area of gracefully degrading systems has been conducted for multiple processor (tightly coupled) systems. These models and the resulting measures, while invaluable to our confidence in tightly coupled multiprocessor systems or loosely coupled systems in which all processors are performing similar functions, prove less meaningful when applied to a large portion of loosely coupled distributed systems such as those used for defense. The main reason is that these analytical models lack consideration for hardware and software topology factors.

Some recent analytical research which addresses the issue of survivability in gracefully degrading systems attempts to accommodate hardware topology and software allocation features of those systems (14,19). One effort by Merwin and Mirhakak proposes that distributed systems are made up of hardware networks and software systems (19). The networks comprise nodes and links. The software system is made up of programs and data. Several failure modes are described. Failure can be caused by loss of a link between a program and its data, loss of the node on which the data resides or loss of the node on which the program is to execute. Failure probabilities are assigned to each node and link. Survivability is determined as the number of programs that remain operational after some combination of nodes and links have failed. A survivability criterion is developed which is based on the

probability of occurrence of any subarchitecture of a given distributed network and the expected number of programs operable for each subarchitecture. A subarchitecture here is defined as any combination of nodes and links which is a subset of the original network configuration. Survivability is expressed as a function of 1.) an initial network architecture, 2.) a given data set distribution and 3.) data set access requirements. The major departure of this work from preceding research is the inclusion of software distribution into the computation of survivability. This model like other analytical models faces several difficulties.

The first problem is computational. In (19) presented above, a number of additions and enhancements to their analytical model are proposed such as weighting of programs and nodes and placing constraints on the data set distribution. However, since the algorithm they use for computing survivability already demonstrates exponential computational growth and complexity as the number of nodes and communications links increase, additional criteria might only serve to exacerbate the present computation problem.

The second problem for the analytical approach is validation. Like other wholly analytical models of distributed systems, the model proposed by Merwin and Mirhakak suffers for lack of validation through fielded systems or experimentation. Although the results are intuitively appealing they are unsubstantiated in application.

A third problem is of particular import to analytical modeling. That is, many system attributes which may be important to distributed system survivability are difficult to measure. Foremost among these

features are those which pertain to software. In earlier fault tolerant systems, software was a minor consideration. Currently software is a primary consideration, and the necessity to incorporate software factors into system evaluation is unavoidable (12).

Whereas some sciences in their early stages are inexact, other sciences are inherently inexact (13). For a philosophical discussion of exact versus inexact sciences, see reference (13). Software is not subject to static standards or metrics but rather must be indirectly described in terms of those attributes which can be measured or observed. Among these attributes are requirements measured in terms of instructions to be executed, storage demands, input/output data rates and application module quantity, size and connectivity. These measurements are used to form the basis for prediction of software related phenomena like development, maintainability and life cycle cost. In that these measurements are rarely derived from first principles, it is unlikely that we can undertake their explanation and prediction in a wholly quantitative manner without imposing severe constraints on the level of complexity to be addressed (7).

Empiricism offers the opportunity to develop statistical laws which can serve several purposes (13). First, it can enhance our understanding of phenomena and provide a basis for prediction or decision making. Second, it can point to areas in which purely analytical investigation can be productive and, third, it can provide a mechanism for validation of analytical models. In addition empiricism does not inherently prejudice research findings and thus establishes essential objectivity.

For these reasons, this dissertation proposes an alternative approach to survivability explanation and estimation. First, a simulation approach to survivability experimentation is taken. Second, a broad spectrum of distributed system attributes are examined. The attributes fall into three general categories namely; those that describe the distributed system capabilities and topology, those that describe the application system topology and requirements and those that describe the distribution and redistribution policies which map the software onto the hardware. In addition, software is permitted to degrade gracefully, that is reduce its resource demands, in order to accommodate degradation of the distributed network. The objective of this dissertation is to provide a foundation for the quantification of operational survivability in gracefully degrading distributed processing systems which can be empirically tested and generally applied.

### CHAPTER III

#### APPROACH

As indicated in the preceding chapter, this research proposes that survivability is a function of a number of attributes. These attributes fall into three general categories, i.e., those that describe the distributed network, those that describe the application system, and those that describe the distribution policy. Further, this dissertation proposes that none of these categories taken alone serve adequately to explain or predict the survivability of a given system. The method of this research is to investigate the relationship between a number of attributes of distributed systems and survival of those systems in the face of increasing degradation of network resources.

Several alternative approaches to this investigation have been considered. The approach must facilitate manipulation of a number of factors pertaining to distributed networks, application systems and distribution approaches for the purpose of analysis. Perhaps the most likely alternative from the point of flexibility is an analytical model. However, the constituents of distributed systems such as routing, resource allocation and task assignment when individually subjected to analytical study present difficult and complex problems. Among these problems are measurement and computability problems. Many system attributes are difficult to describe quantitatively such as reconfiguration options or decisions, resource capabilities and

execution constraints. When quantitative description is possible, it often shows exponential growth with increases in system size. It follows, therefore, that a system comprising many of these constituents would be correspondingly more difficult to represent analytically (19,14). Further, given an analytical approach a decision must be made to either oversimplify the distributed processing problem or address the potential problem of intractability.

The second major alternative is empirical. If an experimental approach is to be used, a choice must be made first between field versus laboratory experimentation. Since instances of operational distributed processing systems are few, opportunities for field experiments at this time are commensurately limited. For these reasons laboratory experimentation was selected as the most viable approach for this research. Laboratory experimentation via simulation permits the representation and control of factors in the field environment which are controllable and many which are not. The degree to which simulation can present a "true" picture of that which it simulates is dependent upon the level of understanding which exists of the components, actions and interactions of the simulated phenomena. Simulation differs from analytical modeling in that the relationships between the proposed factors need not be stated in an explicitly quantitative fashion. Once a simulation facility exists, the subsequent capability for controlled replication takes major advantage over field experimentation. Also, the ability to monitor and track operational conditions, decisions, and intermediate actions can be considerably easier than in field situations.

Computer simulation has been chosen for examination of the problem of distributed system survivability because of the large number of variables which must be used, the complexity of their manipulation and the magnitude of the possible instantiations of the variable values.

Computer simulation, allows a "systems view" of distributed systems to be presented. The first objective of this research is to determine those factors which have an important effect on distributed system survivability and can be used in the development of a measure or model of survivability. Second, this research proposes to provide a simulation approach to distributed system comparison. Since experimentation via simulation is similar in many ways to field test, traditional experimental methods and analysis techniques should apply. One intent of this research is, in fact, to demonstrate the applicability of standard experimental design and analysis techniques to topics in computer science.

The following provides an introduction and discussion of the experimental approach to be used. Specifics of implementation in the present investigation are given in Chapter IV.

The literature repeatedly paints the picture of distributed systems as a very complex one comprised of numerous orthogonal and interrelated factors such as levels of hardware, control and load base decentralization (8,15). One goal of this work is to discover the subset of active factors which is important to survivability. The process of factor selection is called factor screening. Factor screening must take place with a comprehensive set of active factors,



because omission of an important active factor can introduce consequences such as bias in the analysis and conclusions drawn from the experiment. Inclusion of negligible factors may, on the other hand, be unnecessarily resource consumptive and introduce sufficient noise in the data such that important effects are difficult to recognize.

Factor screening methods can be introduced either during the design and development of the simulation model or upon completion of the simulator. When employed at early stages, as is proposed here, factor screening will affect the choice of variables and variable levels in the model. The overall impact will be to simplify the structure of the final model and sharpen description of specific effects (16).

Let us suppose that there are a number of controllable factors in the simulation experiment, call these  $X_1, X_2, \dots, X_K$  and two response variables  $S$  and  $P$ . Since  $S$ , survivability, for now is assumed to be two-valued and be a function of some range in  $P$ , performance, such that

$$S = \begin{cases} 1 & \text{if } P \leq n \\ 0 & \text{otherwise} \end{cases} \quad (3-1)$$

we can proceed as though there is but one response,  $P$ . Further let us assume that the simulator is structured such that the response can be expressed in the form

$$P = f(X_1, X_2, \dots, X_K) +$$

where  $f$  is a function that determines the mean value of  $P$ , and represents error such that,  $E(\epsilon)=0$ , the expected value of  $\epsilon$  is zero. Initially, it is assumed that  $f$  is linear in the unknown parameters, coefficients, that relate the response,  $P$ , to the factors,  $X_1, X_2, \dots, X_K$ . One possible model is

$$P = B_0 + \sum_{i=1}^K B_i X_i + \epsilon \quad (3-2)$$

where  $B_0, B_1, \dots, B_K$  are unknown parameters. Here  $B_0$  is the intercept and  $B_1, B_2, \dots, B_K$ , the coefficients.

To use this system to conduct an experiment, the levels of each factor must be chosen and the simulation run on the full set or some subset of the factor level combinations. The selection of the number of factor levels to be used and their spacing is extremely important. Since, in this research, as in many factor screening experiments, we are trying to determine the relative effect of a factor and not develop a highly precise predictive or interpolative equation, the number of factor levels or values to be tested will be small, two or four. The "effect" of a factor is described as the change in the response observed as a result of a change in levels of the factor. This direct cause-effect relationship between a single factor and the response is called a "main" effect.

Factor screening experiments fall into two major categories, full factorial experiments and fractional factorial experiments. The most efficient full factorial design is the  $2^K$  factorial design which

comprises  $K$  factors each having two levels. The statistical model generated for a  $2^K$  full factorial design would include  $K$  main effects,  $\binom{K}{2}$  two-factor interactions,  $\binom{K}{3}$  three-factor interactions,  $\binom{K}{4}$  four-factor interactions, etc. and one  $K$ -factor interaction. In total the  $2^K$  design would describe  $2^K - 1$  effects.

The term treatment combination is used to refer to the aggregate of factor settings of all factors as designated for a given experiment run or case. One system of notation frequently used to denote individual factor levels, uses + and - signs to designate high and low or alternate levels of the factor. Thus, a treatment combination for a four factor experiment on factors  $X_1, X_2, X_3, X_4$  might be - + - - indicating that factors  $X_1$  and  $X_4$  are at their low setting and factors  $X_2$  and  $X_3$  at their high setting.

The total number of experiment runs required in a  $2^K$  full factorial design given small values of  $K$  such as 5 or 6 is 32 and 64 respectively. The magnitude of this number grows exponentially with  $K$ . Since resources are usually limited, the number of replicates that the experimenter can employ may be restricted. Frequently, available resources will only allow a single replicate of the design to be run, unless the experimenter is willing to omit some of the original factors.

With only a single replicate of the  $2^K$  it is impossible to compute an estimate of experimental error, that is, a mean square for error. Thus, hypotheses concerning main effects and interactions cannot be tested. However, the usual approach to the analysis of a single replicate of a  $2^K$  full factorial design is to assume that

certain higher-order interactions are negligible (21). The statistical analysis of these designs by either Yates' tabular algorithm or a regression approach may be used to estimate the effects. Since this is a factor screening experiment, our interest will be confined to detecting main effects and 2-factor interactions. We can, therefore, either use the higher-order effects as an estimate of error, or as the basis of developing a more efficient design via fractional replication. By assuming that certain high-order interactions are negligible, information on main effects and low-order interactions may be obtained by running only a fraction of the complete factorial experiment. These fractional factorial designs are widely used in research and have major applications in factor screening (21).

In a  $2^{K-P}$  fractional factorial design, only a fraction,  $1/2^P$ , of the  $2^K$  treatment combinations are actually run. A fraction of the  $2^K$  design containing  $2^{K-P}$  runs is called a  $1/2^P$  fraction of the  $2^K$  full factorial design, or a  $2^{K-P}$  fractional factorial design. The design proposed in this research is a regular fraction, that is, estimates of the effects are orthogonal. The effects may be estimated by generating the contrast for any factor using the table of + and - signs for that design which is equivalent to the regression approach outlined above. There are several commonly used methods of constructing these designs.

The particular  $2^{K-P}$  fractional factorial design to be used in this research is of resolution  $V$  usually expressed as  $2^{K-P}_V$ . In a resolution  $V$  design an unconfounded estimation of all main effects and two factor interactions is obtained. Three factor interactions and higher will be confounded or aliased in such a way that isolation of

particular effects is not possible. The higher the resolution of the fractional factorial design, the greater the information obtained concerning higher order interactions. The higher the resolution, the closer the fractional factorial design comes to a full factorial design and consequently the greater the number of experiment runs required. It follows that as the size of  $K$  increases, the number of experiment runs required to meet higher resolution designs is directly affected. Selection of the appropriate design resolution is an important part of initial research considerations. For further information on  $2^{K-P}$  fractional factorial designs the reader is referred to two papers by Box and Hunter (6,9).

## CHAPTER IV

### PROCEDURE

The rationale for the simulation approach to model design was presented in Chapter III. The objective of this simulation is to facilitate determination of those factors which have an important effect on distributed system survivability and can be used to develop a measure or index of survivability. In addition it is anticipated that this simulation approach will be used to compare distributed processing systems. A discussion of the initial simplifying assumptions, variable selection and quantification is provided below. In addition, the basic  $2^{K-P}$  fractional factorial resolution V experimental design is described, the experimental approach outlined and the basic structure of the simulator is presented.

### Assumptions

To effect a simulation which comprises adequate variables to represent a realistic distributed processing system and sufficiently well specified to permit experimentation, three simplifying assumptions on the distributed system attributes are used. As experimentation with the proposed simulator proceeds it may be possible to relax some of these constraints. The initial assumptions are discussed below.

1. All software support resources and application software is accessible by all processing nodes. It is, of course, not likely that these resources are all equally easy to access;

however, the complexity of accessibility will not be addressed in this research. This assumption is made so that the issue of application and support software transfer from one node to another need not be addressed. This assumption is realistic for application and support software on homogeneous networks but falls short when changing data bases are considered. This assumption as it relates to data bases will be relaxed in future experiments.

2. Loss of communication links alone is not considered in these experiments. Loss of a node will, of course, eliminate all links connecting to that node. The effects of link loss is a very complex problem which continues to be extensively researched in connection with various types of networks (9,10,11,1). The loss of individual links can be readily incorporated into the experiment setting proposed for this research. In essence since the removal of a node implies the removal of all adjoining links, creation of an artificial node representing a given link and the subsequent removal of that node will have the same effect as removal of the original link.

3. The simulator has control over vulnerability. The vulnerability and criticality of individual processing nodes are very important considerations for many applications and can be incorporated in the proposed simulator at a later

date. Presently, however, omitting these factors allows us to focus on the structural features of distributed systems which effect operational survivability. Both static and dynamic vulnerability and criticality attributes will be added in later experiments.

#### Experiment Factors and Factor Levels

A distributed processing system is a computer network composed of two or more autonomous processing and memory components working together to serve a common application. A gracefully degrading system is a multiple processor system which provides a high quality of service by reconfiguring the system or network or by reallocating resources when a fault is detected. Operational survivability, then, is an attribute describing the degree to which a distributed processing system can gracefully degrade. The objectives of this research are to make our understanding of survivability a quantitative one and to develop a model or set of models with which we can evaluate and predict operational survivability and performance. The survivability index can be expressed as a simple function of level of performance. In this research, performance can have one of four values depending on the level to which application system requirements are satisfied.

Performance value

- 1 indicates normal or satisfactory application system performance
- 2,3 indicate satisfactory degraded application performance
- 4 indicates unsatisfactory application system

performance.

Satisfactory degraded application system performance refers to the success of the distributed system to adjust to a loss of distributed network resources by a reduction in application system requirements. The survivability index will have either the value "1" indicating that a given distributed system is survivable or 0 indicating that the system is not survivable according to the following

$$\text{Survivability Index} = \begin{cases} 1 & \text{if Performance} \leq 3 \\ 0 & \text{otherwise} \end{cases} \quad (4-1)$$

The value assigned to performance can be expressed as a function of a number of attributes

$$P = f(Z_1, Z_2, \dots, Z_k) \quad (4-2)$$

such that attributes  $Z_1, Z_2, \dots, Z_k$  describe features of the distributed network, application system and distribution policy. The  $Z$ s represent features of the distributed system which are manipulated or controlled.

The parameters that will be controlled in the proposed  $2^{k-p}_v$  Fractional Factorial design are presented in Table 1.

Table 1. Experimental Factors and Factor Levels

	Factor	Levels
Z1	Type of Distributed Processing Topology	a. STAR b. RING c. NETWORK d. ARRAY
Z2	Number of Nodes	a. 4 b. 10
Z3	Node Processing Speed	a. 500 KOPS b. 10 MOPS
Z4	Node Memory Capacity	a. 128 KBYTES b. 2 MBYTES
Z5	Connectivity of Application System	a. Low b. High
Z6	Number of Application Modules	a. 4 b. 16
Z7	Average Module Processing Requirements	a. 10% Node Processing Capacity b. 50% Node Processing Capacity
Z8	Average Module Memory Requirements	a. .1 Node Memory Capacity b. .8 Node Memory Capacity
Z9	Average Frequency of Module to Module Interaction ( $\Delta$ of thousand message set ups)	a. Low b. High
Z10	Distribution/Redistribution Strategy	a. Random b. Uniform c. Packed d. Optimal Spare
Z11	Percent of Nodes Eliminated	a. 10% b. 30% c. 50% d. 80%

Note: For further description of these factors see Appendix A.

Table 2. below shows the correspondence between the eleven variables in the preceding chart and the pseudo factors used in the  $2^{K-P}$  design proposed here. The pseudo factors are used to create 2 two-level factors to represent each four level factor.

Table 2. Experiment Factors and Pseudo-factors

ORIGINAL FACTORS	NO. LEVELS	PSEUDO-FACTORS	LABELS
$Z_1$	4	$X_1$ $X_2$ *	A B
$Z_2$	2	$X_3$	C
$Z_3$	2	$X_4$	D
$Z_4$	2	$X_5$	E
$Z_5$	2	$X_6$	F
$Z_6$	2	$X_7$	G
$Z_7$	2	$X_8$	H
$Z_8$	2	$X_9$	I
$Z_9$	2	$X_{10}$	J
$Z_{10}$	4	$X_{11}$ *	K
		$X_{12}$	L
		$X_{13}$	M
$Z_{11}$	4	$X_{14}$ *	N O

\* Considered Together

Thus, according to Table 2. it is apparent that factors  $Z_1, Z_{10}$  and  $Z_{11}$  are decomposed to 2 two-level pseudo factors. When designing experiment runs, pairs of pseudo factors are considered together.

Let us consider now the design of a fourteen factor experiment. Since this research is concerned with both main effects and two factor

interactions, the Resolution V design is considered. This design provides the desired clarity of main effects and two factor interactions. Implementation of this design for our fourteen factor experiment proceeds as follows. There are fourteen main effects and  $(2^{14})$  or 91 possible two factor interactions which gives a total of 105 effects. Taking the next higher power of two,  $2^7$  indicates that 128 experimental runs would have to be made to cover all the effects of interest. Thus, rather than  $2^{14}$  runs only  $2^{14-7}$  runs, or 1/128 of the total possible combinations need be tried.

Next, it is necessary to describe the individual runs or treatment combinations which must be executed. To construct the chart of experiment runs shown in Table 3. first the plus and minus levels for a full  $2^7$  design in A, B, C, D, E, F, and G is established. Letters here represent factors. The levels for the 7 remaining factors are generated using interactions of the original seven factors as follows:

H=ABCG, J=BCDE, K=ABDF, L=AEFG,  
M=CDFG, N=ACEFG, and O=BDEFG.

Thus, the generating relations for this design are

I=ABCGH, J=ABCDEJ, I=ABDFK, I=AEFGL,  
I=CDFGM, I=ACEFGN, and BDEFGO.

Table 3.

101251 1949.1.28? 17 x 1.6. 75 (J-22)

WJLV 1.1m/s ALL 1174.647ms

FYRCLWCT B.L. RSRQTPPTLCS

[illegible]

Table 3 continued.

2015350 1818062983 181110965 (307)2

STANDARDIZATION OF THE

2011.11.17 15:30

[illegible]

Table 3 continued.  
 2(4-6) FRACTIONAL FACTORIAL DESIGN  
 SURVIVABILITY SIMULATION  
 EXPERIMENT RUN DESCRIPTIONS

RUN	1	2	3	4	5	6	7	8	9	10	11	12
861	+	+	+	+	+	+	+	+	+	+	+	+
871	+	+	+	+	+	+	+	+	+	+	+	+
881	+	+	+	+	+	+	+	+	+	+	+	+
891	+	+	+	+	+	+	+	+	+	+	+	+
901	+	+	+	+	+	+	+	+	+	+	+	+
911	+	+	+	+	+	+	+	+	+	+	+	+
921	+	+	+	+	+	+	+	+	+	+	+	+
931	+	+	+	+	+	+	+	+	+	+	+	+
941	+	+	+	+	+	+	+	+	+	+	+	+
951	+	+	+	+	+	+	+	+	+	+	+	+
961	+	+	+	+	+	+	+	+	+	+	+	+
971	+	+	+	+	+	+	+	+	+	+	+	+
981	+	+	+	+	+	+	+	+	+	+	+	+
991	+	+	+	+	+	+	+	+	+	+	+	+
1001	+	+	+	+	+	+	+	+	+	+	+	+
1011	+	+	+	+	+	+	+	+	+	+	+	+
1021	+	+	+	+	+	+	+	+	+	+	+	+
1031	+	+	+	+	+	+	+	+	+	+	+	+
1041	+	+	+	+	+	+	+	+	+	+	+	+
1051	+	+	+	+	+	+	+	+	+	+	+	+
1061	+	+	+	+	+	+	+	+	+	+	+	+
1071	+	+	+	+	+	+	+	+	+	+	+	+
1081	+	+	+	+	+	+	+	+	+	+	+	+
1091	+	+	+	+	+	+	+	+	+	+	+	+
1101	+	+	+	+	+	+	+	+	+	+	+	+
1111	+	+	+	+	+	+	+	+	+	+	+	+
1121	+	+	+	+	+	+	+	+	+	+	+	+
1131	+	+	+	+	+	+	+	+	+	+	+	+
1141	+	+	+	+	+	+	+	+	+	+	+	+
1151	+	+	+	+	+	+	+	+	+	+	+	+
1161	+	+	+	+	+	+	+	+	+	+	+	+
1171	+	+	+	+	+	+	+	+	+	+	+	+
1181	+	+	+	+	+	+	+	+	+	+	+	+
1191	+	+	+	+	+	+	+	+	+	+	+	+
1201	+	+	+	+	+	+	+	+	+	+	+	+
1211	+	+	+	+	+	+	+	+	+	+	+	+
1221	+	+	+	+	+	+	+	+	+	+	+	+
1231	+	+	+	+	+	+	+	+	+	+	+	+
1241	+	+	+	+	+	+	+	+	+	+	+	+
1251	+	+	+	+	+	+	+	+	+	+	+	+
1261	+	+	+	+	+	+	+	+	+	+	+	+
1271	+	+	+	+	+	+	+	+	+	+	+	+
1281	+	+	+	+	+	+	+	+	+	+	+	+
1291	+	+	+	+	+	+	+	+	+	+	+	+
1301	+	+	+	+	+	+	+	+	+	+	+	+

To determine the level for each 2 level factor simply interpret the corresponding plus or minus sign. To determine the level for four level factors the following set of index tables will be used.

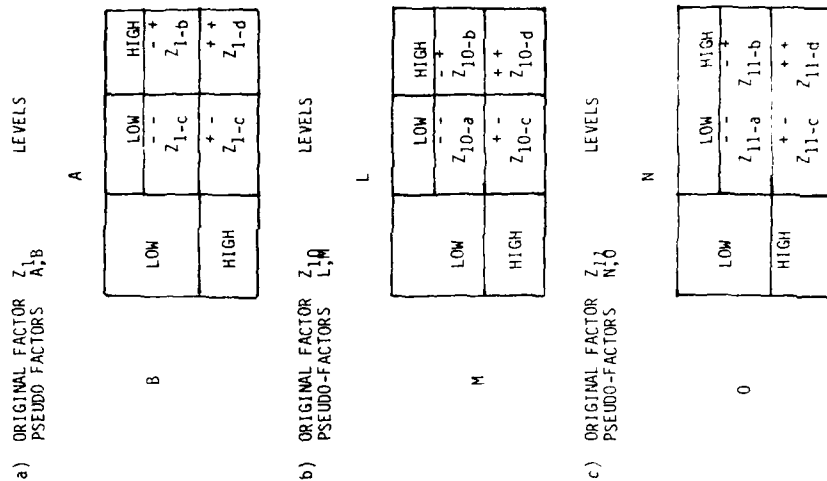


Figure 1. Index Tables for Four Level Factors



Using Tables 1,2 and 3 and Figure 1, run #1 of this experiment would be composed as follows:

Table 4. Interpretation of Example Treatment Combination

RUN #1	
Z <sub>1</sub>	STAR
Z <sub>2</sub>	4 NODES
Z <sub>3</sub>	500 KOPS PROCESSING
Z <sub>4</sub>	128 KB ITS MEMORY CAPACITY
Z <sub>5</sub>	LOW APPLICATION SYSTEM CONNECTIVITY
Z <sub>6</sub>	LOW # APPLICATION MODULES
Z <sub>7</sub>	50,000 KOPS/EXECUTION AVERAGE MODULE PROCESSING REQUIREMENTS
Z <sub>8</sub>	AVERAGE MODULE USES .8 OF NODE MEMORY CAPACITY
Z <sub>9</sub>	HIGH # OF MESSAGE SET UPS
Z <sub>10</sub>	OPTIMAL SPARE DISTRIBUTION
Z <sub>11</sub>	10% NODES ELIMINATED

# SURSIM Survivability Simulator

SURSIM is a simulator which facilitates the investigation of the concept of survivability in gracefully degrading systems. It examines distributed system resources, processing nodes and associated links, which can be lost before a given application system required to execute on that distributed system must function in a degraded mode or experience failure.

The Survivability Simulator depicted in Figure 2 shows the function and flow of the system. SURSIM accepts the description of arbitrary application system topologies and requirements, and distributed system topologies and capabilities, and using predetermined configuration and reconfiguration strategies exercises the hardware/software systems through a sequence of hits or node losses which reduce the capability of the distributed processing system. Effects of configuration modification and capability reduction on application system performance is analyzed. Based on this analysis the application system is reconfigured or the distributed system is further mutated. The simulator continues to iterate reconfigurations and mutations while logging performance and configuration data until the distributed system fails, i.e. the application system can no longer function on the distributed system at an acceptable level.

Within the simulator, the application system and distributed processing network are represented as graphs. For the application system the vertices represent program modules and the edges represent module interaction. For the distributed processing network the vertices represent processing nodes and the edges represent

communication links. Application system requirements are described in terms of module memory requirements, processing requirements, frequency of execution, frequency of module to module interaction, and module criticality. The capability to systematically reduce application system demands according to some a priori defined policy exists. The degree to which procedures of the application system degradation policy are implemented depends upon the degradation level of the distributed network. Distributed system capabilities are described in terms of node processing speed, memory size, and communications capacity. Several different approaches to task assignment are simulated. Via these policies the application system is mapped onto the distributed processing network. This is a graph mapping which is performed according to one of four policies. The four policies are 1.) random distribution, 2.) uniform distribution, 3.) packed distribution and 4.) the optimal-spare distribution. In the random distribution, application system modules are randomly assigned to processors. This will be repeated until all modules have been assigned to nodes or the policy fails to construct a map. If the application module and communication burden exceed that of the node selected, assignment will not be made. In the uniform distribution, application system modules are assigned to nodes such that each node has as near the same operating demands as possible. In the packed distribution, application system modules are assigned to a designated processor until it reaches maximum capacity after which modules are assigned to the "next" processor, etc. In the optimal-spare distribution, application system modules are assigned to the distributed processing system explicitly by

the system designer. Each node being assigned application tasks has a spare queue indicating the sequence of backup or spare nodes which will be activated should the former fail. This distribution approach takes into account the requirement of certain application modules for processing nodes with special I/O devices such as sensors and actuators.

The performance analyzer performs a comparison of application system requirements to the specific distributed system capabilities assigned to it. For each node in the distributed system a comparison is made between the node capability and the application system requirements of all modules assigned to it. For example, if the memory capacity of a node less the memory requirements of all modules assigned to it gives a negative result, performance is considered unsatisfactory. Likewise, if the processor demands exceed the processor capability performance is considered unacceptable. The ability of communications links to meet expected demands is similarly determined by accessing resource saturation. Should the performance analyzer indicate that performance in the current application system/distributed system configuration is satisfactory in all categories, the distributed system topology will be further mutated, otherwise the application system reconfiguration segment of the simulator will be instantiated.

The function of the distributed system topology mutator is to systematically eliminate nodes and their associated links until the distributed system topology is such that "satisfactory" application system performance cannot be achieved. The approach is as follows. First, each individual node and its associated links is removed, then

all possible combinations of two nodes, then three node combinations, etc., until all possible mutations of the distributed system topology have been exercised. The loss of multiple nodes is treated as though these losses occur simultaneously; however, a more advanced form of the simulator should be able to take into account history dependence of failures.

The function of the application system reconfiguration segment of the simulator is to carry out the distribution policy in effect and institute the degradation procedures as necessary. This simulator segment is called into operation when the application system performance analyzer indicates an unsatisfactory level of application system performance. An attempt will be made to reconfigure the application system using whatever distribution policy is in effect to bring the system to an acceptable performance level. Should the reassignment efforts fail to bring performance to the desired level the a priori stated procedures for software degradation will be imposed. Following instantiation of each degradation procedure, performance will be reevaluated. This process is iterated until satisfactory degraded performance is achieved or all degradation procedures have been implemented. In the latter case, the distributed system will have failed to meet the application system performance requirements in normal or degraded mode and consequently, will be considered inoperable.

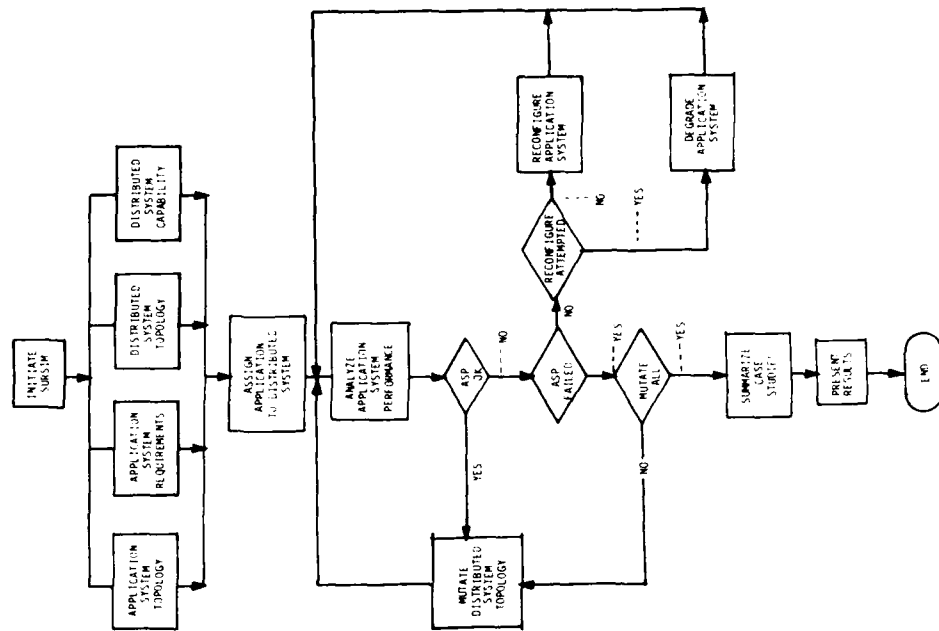
The performance analysis routine in part determines the class of service being provided. Two categories of acceptable service may exist: normal or degraded, and one category of unacceptable service:

failed. There are two ways in which the application system can go from normal to degraded mode. One is to "degrade" or reduce performance requirements. This essentially means that the application modules will continue to perform all their current functions but at a slower rate. The other means of degradation is to "cripple" the application system or purge designated application modules. To what level processing or interaction requirements are reduced or which modules are purged and in what order is determined a priori by the application system designer. This information is input to the simulator. The distributed system continues to be systematically changed by the distributed system topology mutator until application system performance has degenerated to an unacceptable level or all mutations of the distributed network have been exercised.

Information collected by the simulator falls into three categories; 1.) status of controlled factors, 2.) status of indirectly controlled factors, and 3.) derived measures. Controlled factors include; type of distributed system topology, size of network or number of nodes; node processing speed, memory and communications capacity; application system size, connectivity and interaction, processing, and memory requirements; distribution strategy and extent of distributed system degradation (mutation). Indirectly controlled factors include connectivity of the distributed system topology, global resource capacity (processing, memory, communications) and available resource capabilities (processing, memory, and communications). Derived information includes a variety of resources: requirements ratios and consistency measures.

SURSIM, the survivability simulator, has been implemented in FLECS, FORTRAN Language with Extended Control Structures, on a Digital Equipment VAX 11/780 and is now being used as a tool for experimentation on a variety of distributed systems.

Figure 2. Survivability Simulator Flow Diagram



## CHAPTER V

## SIMULATOR RESULTS

Output generated by the simulator for the 128 designed experiment runs and 300,000 subcases fall into two categories. The first type of output is strictly descriptive of the cases run. The second type of output is a log of operational data collected for each of the cases and subcases. Discussion and examples of the output follow.

Descriptive Output

Table 5 is a chart generated by the simulator for each of the 128 designed experiment runs. It presents a string of + and - signs which designate the factor level of each of the original 11 factors for that case followed by an English language interpretation of the factor levels. Table 6 is a representation of the application system inter-actions for this case. Table 7 lists the respective application system requirements for each application module. Table 8 presents the degradation procedures to be followed in the event that the application system cannot function in a satisfactory manner on the distributed system network given its current configuration and application system assignments. Table 9 presents an interaction incidence matrix which describes the topology of the distributed network for this case. The queue of start nodes lists all the potentially unique start nodes for the topology being studied. Table 10 represents the capability of each node in the distributed system and its connection to other nodes. Tables 11 and 12 describe the remaining resources after initial

assignment of the application system to the distributed system and module to node assignments respectively.

Operational Data Logged

Output of the simulator which is descriptive of control factors or records operational data is logged for later analysis. Examples of this output are shown in Table 13. Definition of logged data items follows:

R <sub>1</sub>	S	- Response variable - survivability 1 = survival 2 = failure
R <sub>2</sub>	P	- Response variable - performance 1 = normal satisfactory performance 2,3 = degraded satisfactory performance 4 = unsatisfactory or failed performance
Z <sub>1</sub>	DIS	- Distributed system topology
Z <sub>2</sub>	N	- Number of nodes in the distributed network
Z <sub>3</sub>	NPS	- Node processing speed
Z <sub>4</sub>	NMC	- Node memory capacity
Z <sub>5</sub>	C	- Node communications capacity
Z <sub>6</sub>	M	- Number of application system modules
Z <sub>7</sub>	MPR	- Average application module processing requirements
Z <sub>8</sub>	NMR	- Average application module memory requirements
Z <sub>9</sub>	MIF	- Average module to module interaction frequency
Z <sub>10</sub>	POL	- Distribution policy
Z <sub>11</sub>	PCT	- Percent nodes lost
Z <sub>12</sub>	SN	- Start node



dispersion is one. That is every application module resides on a different node. The closer this ratio comes to 1/m the less dispersed the application system is said to be.

CDISPR,  $Z_{38}$ , has a similar interpretation, however, this measurement is taken at the end of the test case or after a certain percent of the nodes are lost.

MCRIIT,  $Z_{39}$ , represents node criticality. This criticality is determined by summing the connectivity of the application system modules on the nodes which are lost and dividing this sum by the total application system connectivity. As this ratio approaches one, the proportion of the application system to be reallocated is increasing. Also, the character of the portion of the application system to be reallocated is described in terms of its need for cohesion.

The consistency measurements  $Z_{47}$ ,  $Z_{48}$  describe the system in terms of memory and processing demands versus unit node memory and processing capacity. A ratio of one or less indicates that all memory or processing demands can be satisfied by a single node. Ratios greater than one indicate the number of nodes necessary to meet the demands. Note no consideration is made here concerning the capability of the system to make distributions which would use resources optimally.

Link consistency,  $Z_{49}$ , varies slightly from the previous two consistency measures in that it relates average module to module interaction frequency to communication link capacity. This ratio indicates what portion of a link's capacity is consumed by average module to module interaction. The closer this ratio comes to one the

more likely modules will have to reside on the same node or have dedicated links.

Other data values derived were generated via transformations during data analysis. These values represent interactions among other variables. Six new variables of this type were created. These values are calculated by multiplication of the values of variables for which interaction is to be determined. They are

$$\begin{aligned}
 Z_{52} &= \text{Interaction among topologies} \\
 Z_{53} &= Z_2 \times Z_3 \times Z_4 \\
 &\quad (\text{Interaction between number of nodes and node} \\
 &\quad \text{processing speed and node memory capacity}) \\
 Z_{54} &= Z_6 \times Z_7 \times Z_8 \times Z_9 \\
 &\quad (\text{Interaction between number of application system} \\
 &\quad \text{modules and average module processing} \\
 &\quad \text{requirements and average module memory} \\
 &\quad \text{requirements and average module to module} \\
 &\quad \text{interaction frequency}) \\
 Z_{55} &= Z_2 \times Z_{18} \\
 &\quad (\text{Interaction between number of nodes and number} \\
 &\quad \text{of lost nodes}) \\
 Z_{56} &= Z_{45} \times Z_{38} \\
 &\quad (\text{Interaction between application system} \\
 &\quad \text{connectivity and dispersion at the end of a} \\
 &\quad \text{subcase}) \\
 Z_{57} &= Z_{45} \times Z_{38} \times Z_{40} \\
 &\quad (\text{Interaction between application system}
 \end{aligned}$$

Table 5. SURVIVABILITY SIMULATOR										
EXPERIMENT RUN DESCRIPTION										
** CASE NUMBER 33 **										
ORIGINAL Z-FACTORS										
RUN	1	2	3	4	5	6	7	8	9	10
										11
-----										
TYPE OF DISTRIBUTION SYSTEM TOPOLOGY:										
										STAP
										4
NODE PROCESSING SPEED:										
										500 KOPS
										128 KBYTES
CONNECTIVITY OF APPLICATION SYSTEM:										
										HIGH
NUMBER OF APPLICATION MODULES FOR										
										4
A GIVEN SIZE PROGRAM:										
AVERAGE MODULE PROCESSING REQUIREMENTS:										
										50% OF NODE PROCESSING SPEED
AVERAGE MODULE MEMORY REQUIREMENTS:										
										80.0% NODE MEMORY CAPACITY
AVERAGE FREQUENCY OF MODULE TO MODULE										
										LOW
INTERACTION (% OF MESSAGE RETURNS):										
DISTRIBUTION/REDISTRIBUTION POLICY:										
										RANDOM
NUMBER OF NODES ELIMINATED:										
										80.0%

connectivity and dispersion at end of subcase and  
distributed network connectivity)



Table 6. Application System Topology  
Interaction Incidence Matrix

	A	B	C	D
A	0.00	0.63	0.13	0.05
B	0.18	0.00	0.19	0.37
C	0.20	0.41	0.00	0.15
D	0.01	0.42	0.74	0.00

Table 7. Application System Requirements

MODULE IDENTIFIER	MEMORY K BYTES	KOPS/ EXECUTION	EXECUTIONS /T	CRITICALITY
A	92.	100.	4.	2.
B	52.	188.	2.	3.
C	91.	139.	1.	1.
D	24.	63.	4.	4.

Table 8. Degradation Policy

STEP	
1	* Degrade modules of criticality equal 1 to .5 CPU executions memory communications
2	* Purge module D
3	* Degrade Modules of criticality less than 3 to .5 CPU executions memory communications
4	* Failure

Table 9. DISTRIBUTED SYSTEM TOPOLOGY INTERACTION INCIDENCE MATRIX						
	1	2	3	4		
1	100	100	0	0		
2	100	100	100	0		
3	0	100	100	0		
4	0	100	0	100		
QUEUE OF STARTNODES: 1, 2,						
TABLE 10. DISTRIBUTED SYSTEM CAPABILITY						
NODE	MEMORY K BYTES	CPU KOPS	# LINKS IN	CAPACITY IN	# LINKS OUT	CAPACITY OUT
1	128	500	1	100	1	100
2	128	500	3	300	3	300
3	128	500	1	100	1	100
4	128	500	1	100	1	100

Table 11. Resources Remaining After Initial Assignment						
Initial Assignment was successful						
Availability Matrix						
NODE	MEMORY K BYTES	CPU KOPS	# LINKS IN	CAPACITY IN	# LINKS OUT	CAPACITY OUT
1	36	100	1	100	1	100
2	76	124	3	300	3	300
3	128	500	1	100	1	100
4	13	109	1	100	1	100
Table 12. Module to Node Assignment						
Modules						
1	2	3	4			
A	A	A	C			
B	B	B	D			

[illegible][illegible][illegible]

Table 13. Sample Data Log



this form, selection of variables is very important. When, as in the case of this research, data has been collected on more variables than may be necessary in the model, established variable selection techniques can be employed to assist in deciding the most suitable variable mix to use in the final model. The consequences of poor variable selection fall into two general categories. The final model may be 1.) useless or misleading because an important variable has been omitted or 2.) unused because the inclusion of extraneous variables has caused it to be cumbersome. Since no single approach to variable selection is guaranteed to produce the "best" model, it is common to utilize several variable selection techniques to aid in the model building process. The following three subsections describe the techniques used in this work. These techniques differ in 1.) the criterion used for selection of independent variables, 2.) the amount of analysis and comparison that is done using subgroups of independent variables and 3.) the type of residual analysis performed. Residual here refers to the difference between the observed and model-generated or fitted values of a response variable.

#### Multiple Linear Regression

The Multiple Linear Regression technique for variable selection estimates the multiple linear regression equation using all of the independent variables. The coefficients of the regression model are estimated by least squares. This technique is performed by a computerized statistical package PMDP-Routine PIR. Output from this program includes a variety of standard statistical measures for each of the variables in the model. One of the measures provided is the  $T$

statistic. The availability of this statistic for all candidate variables facilitates use of a variable selection approach called the directed search on  $T$ . The test statistic  $T$  will be large for those regressors which contribute significantly to the full model. If these regressors are introduced to the model one at a time in order of descending  $T$  value, the model should be at any given point the "best" or one of the best for that size subset of all possible regressors. The directed search on  $T$  is a good variable selection strategy when the number of variables is large, say 20 to 30.

#### Stepwise Regression

The stepwise regression technique enters and removes variables from a multiple linear regression equation in a stepwise manner. At each step in the model building process variables are removed and/or entered into the equation. The criteria for determining entry or removal of a variable is normally its  $F$  statistic when considered along with other variables in the model. Forward stepping is an approach which begins with no predictors and consecutively adds variables which exceed some threshold value. Backward stepping is an approach which begins with all candidate predictors and consecutively removes variables which fall below a given lower bound. Techniques used in this research employ a combination of forward selection and backward elimination.

#### All Possible Subsets Regression

The all possible subsets regression procedure requires the fitting of all of the regression equations involving one through  $n$  candidate regressors, where  $n$  is the number of variables. The number

of equations to be examined increases exponentially with the number of candidate regressor variables. To evaluate subset regression equations several measures can be used. These include  $R^2$ , coefficient of multiple determination; adjusted  $R^2$ , minimal residual mean square;  $MS_E$ , mean square for error; and Mallows  $C_p$ , residual sum of squares. All possible subsets regression using adjusted  $R^2$  and Mallows  $C_p$  are used here.

In the adjusted  $R^2$  evaluation, the T statistic for the coefficients of variables in the subset that maximizes adjusted  $R^2$  are all greater than one in absolute value. Maximizing adjusted  $R^2$  is the same as minimizing the residual mean square. Usually, subsets larger than those that maximize adjusted  $R^2$  are not very good. In the Mallows  $C_p$  evaluation, the T statistic for the coefficients of variables in the subset that minimizes  $C_p$  are usually greater than  $\sqrt{2}$  in absolute value. When using all possible subsets regression the problems of variable selection increase as the number of redundant variables increases. Inclusion of irrelevant variables provides the opportunity for artifacts in the data to produce unpredictably high T statistics,  $R^2$ , and adjusted  $R^2$  and unpredictably low Mallows  $C_p$  statistics. For this reason checks must be made for variable redundancy and redundant variables removed from the set of candidate variables.

#### Procedures for Model Building

##### Data Reduction

The experimental design presented in Chapter IV described the 128 experiment runs necessary for a  $2^{k-p}$  design in 14 factors. During execution of the simulator data was collected for each of these cases

plus two types of subcases. The subcases were those that tracked operational data for all possible unique start nodes and all possible number of nodes lost. The total number of cases and subcases logged by the simulator are in excess of 300,000. Physically this translates to approximately 90 megabytes of data which is one very large magnetic disk or seven 2,400 foot 1,600 BPI magnetic tapes. The mechanical difficulty of working with this volume of data suggests that the possibility of meaningful data reduction should be explored. Fortunately, some reduction of the data could be performed without significantly decreasing its information value. Therefore, before analysis the raw data was put through a data reduction filter which produced three sets of data, each of different resolution.

DATA 1 - comprises the 128 designed experiment runs times an averaging over all possible start nodes times an averaging over number of nodes lost. The size of this data set is 2,156 cases.

DATA 2 - comprises the 128 designed experiment runs times an averaging over number of nodes lost. The size of this data set is 715 cases.

DATA 3 - comprises the 128 designed experiment runs. The size of this data set is 128 cases.

DATA 1 has, of course, the highest resolution of the 3 data sets. It presents a summary of individual subcases such that specific information is lost concerning individual subcases for each start node

distribution policy. Both of these variables have four levels and thus require three "dummy" variables to represent them. This is accomplished by arbitrarily assigning one of the following codes to each of the qualitative variables.

DISTRIBUTED SYSTEM TOPOLOGY	"DUMMY" VARIABLE		
	1A	1B	1C
STAR	0	0	0
RING	1	0	0
NETWORK	0	1	0
ARRAY	0	0	1

DISTRIBUTION POLICY	"DUMMY" VARIABLE		
	1D	1E	1F
RANDOM	0	0	0
UNIFORM	1	0	0
PACKED	0	1	0
OPTIMAL SPARE	0	0	1

The interpretation given to coefficients of qualitative variables is different than that of quantitative variables in that the coefficient of a qualitative variable indicates the relative impact of change from that level to other possible levels of the qualitative variable. For example, in the case of topology each of the "dummy" variables when present in the model indicate the effect of change from the base level or condition, 000, to that level. The effect of change from one of the other levels to a third level is accomplished by subtracting the coefficients of the variables in question. Thus the effect of a change

and each possible number of nodes lost. DATA 2 presents a summary which ignores the start node and DATA 3 presents a summary which ignores both start node and specific number of lost nodes. DATA 3 refers to node loss as a percent of the total nodes initially in the distributed system. Examination of the analyses conducted using all three data sets revealed that the designed data set, DATA 3, was representative of the other two.

#### Candidate Variables

The variables submitted to data analysis are of three types: response variables, control variables, and other independent variables. The response variables are survivability, S, and performance, P. Since S can be obtained from a simple function on P, the focus of discussion of analyses performed will be on P. The control variables are the 11 factors described in Chapter IV section 2. Other independent variables or potential control variables are those listed in Chapter V section 2.

Independent variables fall into two general categories: quantitative or continuous valued variables and indicator variables. Most often variables used in regression model building are quantitative or continuous valued variables which take on values within some known range on a well-defined scale. Less frequently, it is necessary to include qualitative variables which have no natural scale of measurement in the regression model. Qualitative variables, often represented as indicator or "dummy" variables are assigned a set of levels to account for the effect that the variable may have on the response. In this research all independent variables are quantitative with the exception of two. These are distributed system topology and

DISTRIBUTED SYSTEM TOPOLOGY	"DUMMY" VARIABLE		
	IA	IB	IC
STAR	0	0	0
RING	1	0	0
NETWORK	0	1	0
ARRAY	0	0	1

from the star to the ring topology is provided by the coefficient on IA, the star to the network by the coefficient on IB, etc. The effect of a change from the ring to the network is determined by subtracting the coefficient of IB from that of IA. The same approach is used for all comparisons. In the case of quantitative variables, the coefficients indicate the direction and magnitude of the relationship between the independent variable and the response.

An important part of regression analysis is variable selection. In the case of distributed processing systems the most appropriate set of regressor variables is not known and little prior experience exists which might help point the way to initial selection. In such cases it is desirable to begin with the most comprehensive set of candidate variables and reduce this number through iterative selection of regressor sets which are "best" according to one of the evaluation criteria listed above.

#### Explanatory versus Predictive Models

Usually, regression models are valid only over the range of the regressor variables contained in the observed data. Over this interval, the regression equation developed may provide a reasonable approximation of the true functional relationship. However, care

should be exercised to assure that the application of a regression model does not exceed its capability. For example, while some regression models may adequately summarize or describe the data from which they were constructed, they may be less serviceable in describing new data. A model which describes the data to which it was fit is called an explanatory model. Measurements can be made which indicate the adequacy of an explanatory model in fitting its data. Checking for explanatory model adequacy can be done via residual analysis, testing for lack of fit, searching for high-leverage or overly-influential observations and a variety of internal consistency checks (20). It should not be assumed that a model which is proved to fit existing data will also be a good predictor for future data. Further, the model that provides the best fit to existing data may not be equally successful in the final application, that is be a successful predictor. To determine how well the explanatory model will serve as a predictor requires that we validate the model. A number of techniques are available for model validation. These include comparison with other results, collection and comparison with new data, and data splitting. The approaches used to develop explanatory and prediction models for operational survivability are presented in the following sections.

#### The Explanatory Model Building Process

Building a regression model is generally an iterative process requiring repeated analyses as improvements in the model structure or additional special features of the data are discovered. Digital computers and established statistical software can be invaluable model building tools. In this case several regression routines comprised in



The BMDP statistical software package are used. These are PIR; Multiple Linear Regression; P2R, Stepwise Regression; and P9R, All Possible Subset Regression.

Initially multiple linear regression is performed using all candidate independent variables. A check is made for multicollinearity among the independent variables. Redundant variables identified by this check are removed from the list of candidate regressors, and the analysis is repeated. The model resulting from this analysis is shown in Figure 3. Also provided in Table 14 are major statistics such as  $R^2$  and Mean Square for Error for the model and T-statistic, mean and standard deviation for each of the regressor variables.

NOTE: See Table 16  
for Variable Key

TABLE 14. Multiple Linear Regression Model with all Candidate Regressor Variables

Figure 3. Multiple Linear Regression Model with all Candidate Regressor Variables

Next, stepwise is performed. This analysis provides an incremental view of the model as it is being developed. The point at which model building using stepwise regression can be considered complete is at the point in which the  $R^2$  value begins to show only nominal increases and the mean square for error,  $MS_E$ , starts to increase. Figure 4 provides a picture of the regression model at this point. A quick validation can be made at this stage. To perform this validation a directed search on  $T$  for the results of  $PIR$  must be conducted. This search constitutes a ranking of candidate regressor variables according to descending values of  $T$ . When this list is compared to the list of regressor variables proposed as a result of stepwise regression analysis, the variables with the largest  $T$  statistic after the direct search on  $T$  should roughly correspond to the variables remaining in the model after stepwise analysis.

Table 14. Statistics from BMDP Multiple Linear Regression Analysis

Variable	Mean	Standard Deviation	St. Dev.	Mean	Minimum	Maximum
X1	.25000	.43471	1.73886	0.00000	0.00000	1.00000
X2	.25000	.43471	1.73886	0.00000	0.00000	1.00000
X3	.25000	.43471	1.73886	0.00000	0.00000	1.00000
X4	7.00000	3.01179	.43026	4.00000	4.00000	10.00000
X5	250.00000	4768.66412	.90832	500.00000	500.00000	10000.00000
X6	1.50000	.50196	.33464	1.00000	1.00000	2.00000
X7	10.00000	6.02358	.60236	4.00000	4.00000	16.00000
X8	45.00000	35.13753	.78083	10.00000	10.00000	80.00000
X9	1.50000	.50196	.33464	1.00000	1.00000	2.00000
X10	.25000	.43471	1.73886	0.00000	0.00000	1.00000
X11	.25000	.43471	1.73886	0.00000	0.00000	1.00000
X12	.25000	.43471	1.73886	0.00000	0.00000	1.00000
X13	42.50000	25.96181	.61087	10.00000	10.00000	80.00000
X14	1.43750	.49803	.34645	1.00000	1.00000	2.00000
X15	7448.00000	7841.09470	1.05278	512.00000	512.00000	20000.00000
X16	27896.63281	33402.45299	1.20059	331.00000	331.00000	97425.00000
X17	5827.65625	5605.86220	1.13354	137.00000	137.00000	19304.00000
X18	2217.51047	1125.45843	.50753	0.00000	0.00000	3800.00000
X19	.46125	.22541	.48869	.20000	.20000	.83000
X20	.35594	.32141	.90300	.04000	.04000	1.05000
X21	.09312	.17906	1.92284	0.00000	0.00000	1.05000
X22	.44250	.35695	.80667	.10000	.10000	1.00000
X23	.59797	.32887	.54997	.06000	.06000	1.00000
X24	4.50000	4.91310	1.09180	.40000	.40000	12.80000
X25	3.00000	2.95776	.98592	.40000	.40000	8.00000
X26	1997.70312	4095.96502	2.05034	0.00000	0.00000	17172.00000
X27	704.59961	1007.87230	1.43042	0.00000	0.00000	3600.00000
X28	.29836	.37377	1.25273	0.00000	0.00000	1.00000
X29	.21992	.31591	1.43648	0.00000	0.00000	1.00000
X30	20250.00000	32163.95008	1.58834	400.00000	400.00000	128000.00000
X31	.20459	.31684	1.54874	0.00000	0.00000	1.00000
X32	.08424	.14263	1.69324	0.00000	0.00000	.67000

NOTE: See Table 16 for Variable Key

The two analyses conducted at this point should serve to reduce the number of candidate regressor variables. All possible subsets regression (PGR) is now performed using the remaining variables. This model building technique is executed first using Mallows Cp as the variable selection criterion, then using adjusted  $R^2$  as the variable selection criterion. Each of these provides detailed information on the five subset models determined to be "best" according to the evaluation criterion in effect and the one model considered optimum. The five models developed using Mallows Cp and adjusted  $R^2$  as evaluation criteria are contained in Appendix B. The two optimum models are presented in Figures 5 and 6 respectively. Table 15 compares the models developed by multiple linear regression, stepwise regression, and all possible subsets regression according to the evaluation data available.

Figure 5. Optimum Model According to Adjusted R Criterion

Variable	Parameter	Estimate	Standard Error	t-Statistic	Probability >  t	Partial R	Partial R Squared	Adjusted R Squared
1	Intercept	1.0000	0.0000	1.0000	0.0000	0.0000	0.0000	0.0000
2	Variable 2	0.1234	0.0123	10.0000	0.0000	0.1234	0.0123	0.0123
3	Variable 3	0.4567	0.0234	19.5000	0.0000	0.4567	0.0234	0.0234
4	Variable 4	0.7890	0.0345	22.8750	0.0000	0.7890	0.0345	0.0345
5	Variable 5	0.2345	0.0456	5.1438	0.0000	0.2345	0.0456	0.0456
6	Variable 6	0.5678	0.0567	10.0000	0.0000	0.5678	0.0567	0.0567
7	Variable 7	0.8901	0.0678	13.1250	0.0000	0.8901	0.0678	0.0678
8	Variable 8	0.3456	0.0789	4.3750	0.0000	0.3456	0.0789	0.0789
9	Variable 9	0.6789	0.0890	7.6190	0.0000	0.6789	0.0890	0.0890
10	Variable 10	0.9012	0.0901	10.0000	0.0000	0.9012	0.0901	0.0901
11	Variable 11	0.1234	0.1012	1.2188	0.2234	0.1234	0.1012	0.1012
12	Variable 12	0.4567	0.1123	4.0625	0.0000	0.4567	0.1123	0.1123
13	Variable 13	0.7890	0.1234	6.3938	0.0000	0.7890	0.1234	0.1234
14	Variable 14	0.2345	0.1345	1.7438	0.0854	0.2345	0.1345	0.1345
15	Variable 15	0.5678	0.1456	3.8938	0.0000	0.5678	0.1456	0.1456
16	Variable 16	0.8901	0.1567	5.6875	0.0000	0.8901	0.1567	0.1567
17	Variable 17	0.3456	0.1678	2.0625	0.0438	0.3456	0.1678	0.1678
18	Variable 18	0.6789	0.1789	3.7938	0.0000	0.6789	0.1789	0.1789
19	Variable 19	0.9012	0.1890	4.7688	0.0000	0.9012	0.1890	0.1890
20	Variable 20	0.1234	0.1901	0.6479	0.5174	0.1234	0.1901	0.1901
21	Variable 21	0.4567	0.2012	2.2719	0.0281	0.4567	0.2012	0.2012
22	Variable 22	0.7890	0.2123	3.7188	0.0000	0.7890	0.2123	0.2123
23	Variable 23	0.2345	0.2234	1.0479	0.2964	0.2345	0.2234	0.2234
24	Variable 24	0.5678	0.2345	2.4219	0.0174	0.5678	0.2345	0.2345
25	Variable 25	0.8901	0.2456	3.6250	0.0000	0.8901	0.2456	0.2456
26	Variable 26	0.3456	0.2567	1.3438	0.1804	0.3456	0.2567	0.2567
27	Variable 27	0.6789	0.2678	2.5375	0.0125	0.6789	0.2678	0.2678
28	Variable 28	0.9012	0.2789	3.2125	0.0000	0.9012	0.2789	0.2789
29	Variable 29	0.1234	0.2890	0.4262	0.6704	0.1234	0.2890	0.2890
30	Variable 30	0.4567	0.2901	1.5750	0.1174	0.4567	0.2901	0.2901
31	Variable 31	0.7890	0.3012	2.6188	0.0099	0.7890	0.3012	0.3012
32	Variable 32	0.2345	0.3123	0.7500	0.4545	0.2345	0.3123	0.3123
33	Variable 33	0.5678	0.3234	1.7562	0.0824	0.5678	0.3234	0.3234
34	Variable 34	0.8901	0.3345	2.6625	0.0099	0.8901	0.3345	0.3345
35	Variable 35	0.3456	0.3456	1.0000	0.3174	0.3456	0.3456	0.3456
36	Variable 36	0.6789	0.3567	1.9031	0.0599	0.6789	0.3567	0.3567
37	Variable 37	0.9012	0.3678	2.4500	0.0156	0.9012	0.3678	0.3678
38	Variable 38	0.1234	0.3789	0.3256	0.7434	0.1234	0.3789	0.3789
39	Variable 39	0.4567	0.3890	1.1738	0.2434	0.4567	0.3890	0.3890
40	Variable 40	0.7890	0.3901	2.0219	0.0438	0.7890	0.3901	0.3901
41	Variable 41	0.2345	0.4012	0.5844	0.5574	0.2345	0.4012	0.4012
42	Variable 42	0.5678	0.4123	1.3769	0.1704	0.5678	0.4123	0.4123
43	Variable 43	0.8901	0.4234	2.1031	0.0374	0.8901	0.4234	0.4234
44	Variable 44	0.3456	0.4345	0.7956	0.4264	0.3456	0.4345	0.4345
45	Variable 45	0.6789	0.4456	1.5231	0.1274	0.6789	0.4456	0.4456
46	Variable 46	0.9012	0.4567	1.9750	0.0500	0.9012	0.4567	0.4567
47	Variable 47	0.1234	0.4678	0.2638	0.7904	0.1234	0.4678	0.4678
48	Variable 48	0.4567	0.4789	0.9531	0.3374	0.4567	0.4789	0.4789
49	Variable 49	0.7890	0.4890	1.6125	0.1074	0.7890	0.4890	0.4890
50	Variable 50	0.2345	0.4901	0.4794	0.6304	0.2345	0.4901	0.4901
51	Variable 51	0.5678	0.5012	1.1331	0.2574	0.5678	0.5012	0.5012
52	Variable 52	0.8901	0.5123	1.7375	0.0854	0.8901	0.5123	0.5123
53	Variable 53	0.3456	0.5234	0.6594	0.5074	0.3456	0.5234	0.5234
54	Variable 54	0.6789	0.5345	1.2688	0.2074	0.6789	0.5345	0.5345
55	Variable 55	0.9012	0.5456	1.6500	0.1000	0.9012	0.5456	0.5456
56	Variable 56	0.1234	0.5567	0.2219	0.8234	0.1234	0.5567	0.5567
57	Variable 57	0.4567	0.5678	0.8031	0.4204	0.4567	0.5678	0.5678
58	Variable 58	0.7890	0.5789	1.3625	0.1724	0.7890	0.5789	0.5789
59	Variable 59	0.2345	0.5890	0.3979	0.6904	0.2345	0.5890	0.5890
60	Variable 60	0.5678	0.5901	0.9656	0.3324	0.5678	0.5901	0.5901
61	Variable 61	0.8901	0.6012	1.4812	0.1374	0.8901	0.6012	0.6012
62	Variable 62	0.3456	0.6123	0.5625	0.5774	0.3456	0.6123	0.6123
63	Variable 63	0.6789	0.6234	1.0875	0.2774	0.6789	0.6234	0.6234
64	Variable 64	0.9012	0.6345	1.4188	0.1574	0.9012	0.6345	0.6345
65	Variable 65	0.1234	0.6456	0.1906	0.8504	0.1234	0.6456	0.6456
66	Variable 66	0.4567	0.6567	0.6938	0.4874	0.4567	0.6567	0.6567
67	Variable 67	0.7890	0.6678	1.1688	0.2474	0.7890	0.6678	0.6678
68	Variable 68	0.2345	0.6789	0.3438	0.7304	0.2345	0.6789	0.6789
69	Variable 69	0.5678	0.6890	0.8125	0.4174	0.5678	0.6890	0.6890
70	Variable 70	0.8901	0.6901	1.2750	0.2074	0.8901	0.6901	0.6901
71	Variable 71	0.3456	0.7012	0.4906	0.6234	0.3456	0.7012	0.7012
72	Variable 72	0.6789	0.7123	0.9531	0.3374	0.6789	0.7123	0.7123
73	Variable 73	0.9012	0.7234	1.2500	0.2125	0.9012	0.7234	0.7234
74	Variable 74	0.1234	0.7345	0.1679	0.8704	0.1234	0.7345	0.7345
75	Variable 75	0.4567	0.7456	0.6125	0.5374	0.4567	0.7456	0.7456
76	Variable 76	0.7890	0.7567	1.0438	0.2974	0.7890	0.7567	0.7567
77	Variable 77	0.2345	0.7678	0.3031	0.7604	0.2345	0.7678	0.7678
78	Variable 78	0.5678	0.7789	0.7250	0.4624	0.5678	0.7789	0.7789
79	Variable 79	0.8901	0.7890	1.1250	0.2574	0.8901	0.7890	0.7890
80	Variable 80	0.3456	0.7901	0.4375	0.6574	0.3456	0.7901	0.7901
81	Variable 81	0.6789	0.8012	0.8438	0.4024	0.6789	0.8012	0.8012
82	Variable 82	0.9012	0.8123	1.1094	0.2674	0.9012	0.8123	0.8123
83	Variable 83	0.1234	0.8234	0.1500	0.8804	0.1234	0.8234	0.8234
84	Variable 84	0.4567	0.8345	0.5438	0.5874	0.4567	0.8345	0.8345
85	Variable 85	0.7890	0.8456	0.9375	0.3474	0.7890	0.8456	0.8456
86	Variable 86	0.2345	0.8567	0.2750	0.7804	0.2345	0.8567	0.8567
87	Variable 87	0.5678	0.8678	0.6500	0.5124	0.5678	0.8678	0.8678
88	Variable 88	0.8901	0.8789	1.0188	0.3124	0.8901	0.8789	0.8789
89	Variable 89	0.3456	0.8890	0.3875	0.6974	0.3456	0.8890	0.8890
90	Variable 90	0.6789	0.8901	0.7625	0.4424	0.6789	0.8901	0.8901
91	Variable 91	0.9012	0.9012	1.0000	0.3174	0.9012	0.9012	0.9012
92	Variable 92	0.1234	0.9123	0.1333	0.8904	0.1234	0.9123	0.9123
93	Variable 93	0.4567	0.9234	0.4906	0.6234	0.4567	0.9234	0.9234
94	Variable 94	0.7890	0.9345	0.8438	0.4024	0.7890	0.9345	0.9345
95	Variable 95	0.2345	0.9456	0.2469	0.8104	0.2345	0.9456	0.9456
96	Variable 96	0.5678	0.9567	0.5938	0.5474	0.5678	0.9567	0.9567
97	Variable 97	0.8901	0.9678	0.9188	0.3574	0.8901	0.9678	0.9678
98	Variable 98	0.3456	0.9789	0.3479	0.7274	0.3456	0.9789	0.9789
99	Variable 99	0.6789	0.9890	0.6875	0.4924	0.6789	0.9890	0.9890
100	Variable 100	0.9012	0.9901	0.9188	0.3574	0.9012	0.9901	0.9901

CONSTANT  
TO  
SQUARED

STATISTICS FOR VARIABLE KEY  
1. SEE STEP 10 FOR VARIABLE KEY  
2. SEE STEP 10 FOR VARIABLE KEY  
3. SEE STEP 10 FOR VARIABLE KEY  
4. SEE STEP 10 FOR VARIABLE KEY  
5. SEE STEP 10 FOR VARIABLE KEY  
6. SEE STEP 10 FOR VARIABLE KEY  
7. SEE STEP 10 FOR VARIABLE KEY  
8. SEE STEP 10 FOR VARIABLE KEY  
9. SEE STEP 10 FOR VARIABLE KEY  
10. SEE STEP 10 FOR VARIABLE KEY  
11. SEE STEP 10 FOR VARIABLE KEY  
12. SEE STEP 10 FOR VARIABLE KEY  
13. SEE STEP 10 FOR VARIABLE KEY  
14. SEE STEP 10 FOR VARIABLE KEY  
15. SEE STEP 10 FOR VARIABLE KEY  
16. SEE STEP 10 FOR VARIABLE KEY  
17. SEE STEP 10 FOR VARIABLE KEY  
18. SEE STEP 10 FOR VARIABLE KEY  
19. SEE STEP 10 FOR VARIABLE KEY  
20. SEE STEP 10 FOR VARIABLE KEY  
21. SEE STEP 10 FOR VARIABLE KEY  
22. SEE STEP 10 FOR VARIABLE KEY  
23. SEE STEP 10 FOR VARIABLE KEY  
24. SEE STEP 10 FOR VARIABLE KEY  
25. SEE STEP 10 FOR VARIABLE KEY  
26. SEE STEP 10 FOR VARIABLE KEY  
27. SEE STEP 10 FOR VARIABLE KEY  
28. SEE STEP 10 FOR VARIABLE KEY  
29. SEE STEP 10 FOR VARIABLE KEY  
30. SEE STEP 10 FOR VARIABLE KEY  
31. SEE STEP 10 FOR VARIABLE KEY  
32. SEE STEP 10 FOR VARIABLE KEY  
33. SEE STEP 10 FOR VARIABLE KEY  
34. SEE STEP 10 FOR VARIABLE KEY  
35. SEE STEP 10 FOR VARIABLE KEY  
36. SEE STEP 10 FOR VARIABLE KEY  
37. SEE STEP 10 FOR VARIABLE KEY  
38. SEE STEP 10 FOR VARIABLE KEY  
39. SEE STEP 10 FOR VARIABLE KEY  
40. SEE STEP 10 FOR VARIABLE KEY  
41. SEE STEP 10 FOR VARIABLE KEY  
42. SEE STEP 10 FOR VARIABLE KEY  
43. SEE STEP 10 FOR VARIABLE KEY  
44. SEE STEP 10 FOR VARIABLE KEY  
45. SEE STEP 10 FOR VARIABLE KEY  
46. SEE STEP 10 FOR VARIABLE KEY  
47. SEE STEP 10 FOR VARIABLE KEY  
48. SEE STEP 10 FOR VARIABLE KEY  
49. SEE STEP 10 FOR VARIABLE KEY  
50. SEE STEP 10 FOR VARIABLE KEY  
51. SEE STEP 10 FOR VARIABLE KEY  
52. SEE STEP 10 FOR VARIABLE KEY  
53. SEE STEP 10 FOR VARIABLE KEY  
54. SEE STEP 10 FOR VARIABLE KEY  
55. SEE STEP 10 FOR VARIABLE KEY  
56. SEE STEP 10 FOR VARIABLE KEY  
57. SEE STEP 10 FOR VARIABLE KEY  
58. SEE STEP 10 FOR VARIABLE KEY  
59. SEE STEP 10 FOR VARIABLE KEY  
60. SEE STEP 10 FOR VARIABLE KEY  
61. SEE STEP 10 FOR VARIABLE KEY  
62. SEE STEP 10 FOR VARIABLE KEY  
63. SEE STEP 10 FOR VARIABLE KEY  
64. SEE STEP 10 FOR VARIABLE KEY  
65. SEE STEP 10 FOR VARIABLE KEY  
66. SEE STEP 10 FOR VARIABLE KEY  
67. SEE STEP 10 FOR VARIABLE KEY  
68. SEE STEP 10 FOR VARIABLE KEY  
69. SEE STEP 10 FOR VARIABLE KEY  
70. SEE STEP 10 FOR VARIABLE KEY  
71. SEE STEP 10 FOR VARIABLE KEY  
72. SEE STEP 10 FOR VARIABLE KEY  
73. SEE STEP 10 FOR VARIABLE KEY  
74. SEE STEP 10 FOR VARIABLE KEY  
75. SEE STEP 10 FOR VARIABLE KEY  
76. SEE STEP 10 FOR VARIABLE KEY  
77. SEE STEP 10 FOR VARIABLE KEY  
78. SEE STEP 10 FOR VARIABLE KEY  
79. SEE STEP 10 FOR VARIABLE KEY  
80. SEE STEP 10 FOR VARIABLE KEY  
81. SEE STEP 10 FOR VARIABLE KEY  
82. SEE STEP 10 FOR VARIABLE KEY  
83. SEE STEP 10 FOR VARIABLE KEY  
84. SEE STEP 10 FOR VARIABLE KEY  
85. SEE STEP 10 FOR VARIABLE KEY  
86. SEE STEP 10 FOR VARIABLE KEY  
87. SEE STEP 10 FOR VARIABLE KEY  
88. SEE STEP 10 FOR VARIABLE KEY  
89. SEE STEP 10 FOR VARIABLE KEY  
90. SEE STEP 10 FOR VARIABLE KEY  
91. SEE STEP 10 FOR VARIABLE KEY  
92. SEE STEP 10 FOR VARIABLE KEY  
93. SEE STEP 10 FOR VARIABLE KEY  
94. SEE STEP 10 FOR VARIABLE KEY  
95. SEE STEP 10 FOR VARIABLE KEY  
96. SEE STEP 10 FOR VARIABLE KEY  
97. SEE STEP 10 FOR VARIABLE KEY  
98. SEE STEP 10 FOR VARIABLE KEY  
99. SEE STEP 10 FOR VARIABLE KEY  
100. SEE STEP 10 FOR VARIABLE KEY

Figure 5. Optimum Model According to Adjusted R Criterion

Variable	Parameter	Estimate	Standard Error	t-Statistic	Probability >  t	Partial R	Partial R Squared	Adjusted R Squared
1	Intercept	1.0000	0.0000	1.0000	0.0000	0.0000	0.0000	0.0000
2	Variable 2	0.1234	0.0123	10.0000	0.0000	0.1234	0.0123	0.0123
3	Variable 3	0.4567	0.0234	19.5000	0.0000	0.4567	0.0234	0.0234
4	Variable 4	0.7890	0.0345	22.8750	0.0000	0.7890	0.0345	0.0345
5	Variable 5	0.2345	0.0456	5.1438	0.0000	0.2345	0.0456	0.0456
6	Variable 6	0.5678	0.0567	10.0000	0.0000	0.5678	0.0567	0.0567
7	Variable 7	0.8901	0.0678	13.1250	0.0000	0.8901	0.0678	0.0678
8	Variable 8	0.3456	0.0789	4.3750	0.0000	0.3456	0.0789	0.0789</

### Prediction Model Building Process

There are a number of ways in which regression models can be validated and their value as predictors evaluated. Methods of model validation fall into three general categories. These are:

- 1) analysis of model coefficients and predicted values including comparisons with prior experience, physical theory, other analytical models or simulation results,
- 2) collection of fresh data with which to investigate the models predictive performance,
- 3) data splitting; breaking the original data into groups and using these observations to predict the model's performance as a predictor.

Data splitting, which is the approach taken here, is accomplished by separating available data into two parts, the estimation data and the prediction data. The estimation data is used to build the regression model. The prediction data is then used to study the predictive ability of the model. This technique is also called cross-validation.

Since the experiment conducted for this research is a "designed" experiment, data splitting can be accomplished in a very straightforward manner. One of the factors in the original eleven that will not be included in the final model is used as the determinant for the split. The variable to be used is  $Z_4$ , absolute memory size.

Based on this factor, the data is split into two groups, let us call them DATA A and DATA B. Using DATA B as the estimation data set, new models are fit using only the variables specified for the optimal

Method/Model	Number of Variables	$R^2$	$MS_E$	DF	Adjusted $R^2$	Mallows Cp
Multiple Linear Regression	32	.8069	3500.4	95	-	-
Stepwise Regression	9	.7392	380618.8	118	-	-
All Possible Subsets						
Regression - A	1	.765890		-	.734536	15.03
	2	.769484		-	.736256	15.32
	3	.777987		-	.741325	15.29
	4	.782770		-	.744554	15.02
	5	.782582		-	.744332	15.11
All Possible Subsets						
Regression - B	1	.797205		-	.749952	18.17
	2	.796944		-	.749630	18.30
	3	.796902		-	.749579	18.32
	4	.799129		-	.749896	19.26
	5	.801106		-	.749905	20.32
R-Squared						
$R^2$						
MS <sub>E</sub>						
DF						
Residual						

Table 15. Comparison of Models Constructed by Three Regression Methods

models designated by the all possible subsets regression. A total of 32 variables are possible in the new models. The descriptions for these variables as they are renamed are given in Table 16. A chart showing which variables are used in which models is provided in Table 17. The 10 new models fit using multiple linear regression on DATA B, the estimation set, are presented in Appendix C. Each of these models is used to predict the response values of DATA A, the prediction set. The adequacy of the fitted models as predictors is determined by an  $R^2$  for prediction computed as follows.

$$R^2 \text{ prediction} = 1 - \frac{\sum_{i=1}^N e_i^2}{S_{yy}} \quad (6-3)$$

where  $e = y - \hat{y}$

in which  $y$  is the observed value of the response

$\hat{y}$  is the fitted value of the response

and

$$S_{yy} = \sum_{i=1}^N (y_i - \bar{y})^2 \quad (6-4)$$

is the corrected sum of squares in which  $\bar{y}$  is the mean of the observed responses

The  $R^2$  prediction can occasionally be modified upward in instances when the fitted values of the response exceed the range of the observed response only by a small percent. This modification introduces at the model prediction evaluation phase the same correction that would otherwise be made when using the model as a predictor.

Table 18 is a comparison of the 10 fitted models in terms of their explanatory and predictive performance.

Table 16. Candidate Regressors - Variable Key

Factor Label Code/Variable	Factor Description
IA X1	Dummy Variables Indicating
IB X2	Distributed System Topology
IC X3	
I6 X4	Number of Nodes in the Distributed System
I7 X5	Node Processing Speed
I9 X6	Node Communications Capacity
I10 X7	Number of Application System Modules
I12 X8	Module Memory Requirements
I13 X9	Module to Module Interaction Frequency
ID X10	Dummy Variables Indicating Distribution Policy
IE X11	
IF X12	
I15 X13	Percent Nodes Lost
I17 X14	Initial Assignment Result
R2 X15	Global Memory Capacity
R4 X16	Available Processing Capacity after Initial Assignment
R5 X17	Available Memory Capacity after Initial Assignment
R6 X18	Available Communications Capacity after Initial Assignment
S1 X19	Distributed System Connectivity
S2 X20	Memory Requirements/Useable Memory Capacity
S4 X21	Communications Requirements/Useable Communications Capacity
S6 X22	Application System Connectivity
S7 X23	Dispersion - Initial (Number of nodes over which an application system is distributed/ (Number of application system modules)
S8 X24	Memory Consistency (Number of application system modules)/ (Average number of application system modules that will "fit" on a node - memory wise)

Table 16 continued. Candidate Regressors

Factor Label Code/Variable	Factor Description
S9 X25	Processor Consistency (Number of application system modules)/ (Average number of application system modules that will "fit" on a node-processorwise)
A4 X26	Available Processing Capacity at End of Subcase
A5 X27	Available Communications Capacity at End of Subcase
A6 X28	Dispersion at End of Subcase (Number of nodes over which the application system is distributed)/ (number of application system modules)
A7 X29	Criticality of Lost Nodes (Sum of the connectivity of the application system modules residing on the lost nodes)/ (Application system connectivity)
X3 X30	Interaction between Number of Application System Modules and Module Processing Requirements and Module Memory Requirements
X5 X31	Interaction between Dispersion at End of Subcase and Application System Connectivity
X6 X32	Interaction between Dispersion at End of Subcase and Application System Connectivity and Distributed System Connectivity





Another method of prediction validation is to reverse the roles of the prediction and estimation sets and check the results for consistency. DATA A then becomes the estimation set and DATA B the prediction set. Multiple linear regression is used to fit models for the estimation set. These models are in turn used to predict observed values in DATA B. The explanatory  $R^2$  values for the models built on DATA A are consistently lower than those built on DATA B. In addition when the  $R^2$  for prediction was assessed, only one of the fitted models proved to be a good predictor. That model, number 10, had the largest number of regressors and an  $R^2$  prediction of .65956. The remaining nine models made almost poor predictions. Upon cross examination of the data in the sets DATA A and DATA B after the split, it was determined that the two sets were equivalent in all respects with the exception of three variables. These three variables were indirectly related to the variable which formed the basis for the split. This relationship, thus, caused the high values of these three variables to be in one set and the low value in the other. These three variables  $X_{15}$ ,  $X_{17}$ , and  $X_{26}$  were indirectly related to the response and were present either individually or in groups in most of the models. Fitted models built on the data set with the low values of these variables were for the most part unstable when used as predictors. When the fitted models were built on the data set with the high values the models were stable, however, consistently underpredicted the performance of the data set having the low values. No correction was made for this underprediction because the adjustment would be unique to predicting into DATA A. It is believed that a more representative

prediction set would reveal a stronger prediction capability than is indicated here. A split of the data such that DATA A and DATA B are equivalent on all variables may be possible and should substantiate further the findings presented here.

## CHAPTER VII

## ANALYSIS PART II - INTERPRETATION

Discussion of Explanatory Prediction and Models

In the 10 best subset models resulting from all possible subset regression analysis, a total of 32 candidate regressor variables are possible. The variables included in each of these subset models is presented in Table 17 and a description of each of the variables is provided in Table 16. Certain variables are found in all 10 models. These are  $X_4$ ,  $X_8$ ,  $X_9$ ,  $X_{10}$ ,  $X_{11}$ ,  $X_{12}$ ,  $X_{13}$ ,  $X_{14}$ ,  $X_{19}$ ,  $X_{26}$ , and  $X_{30}$ , which represent number of nodes in the distributed system, module memory requirements, module to module interaction frequency, distribution policy, percent nodes lost, initial assignment result, distributed system connectivity available processing capacity at the end of the subcase and the interaction of all application system related variables. Table 19 presents the coefficients for the variables in the 10 best subset models. As can be observed in this table the coefficients for the nine variables found in all 10 models are approximately equivalent in sign and magnitude for all models. In fact, there exists extreme stability of all coefficients across models. Changes when they occur are proportional. Equivalent signs and magnitudes means that the regression coefficients are good estimates of the effects of these factors upon performance. Also, these variables form the core of a model that will likely be good for both explanatory and predictive assessment. In other words, the 32 variables included

in these models are very stable and are not distorted much by the introduction or removal of other variables.

The number of variables in addition to the nine foundation variables needed to achieve explanatory models with  $R^2$  adequacy levels above .8 vary between four and 11. It is important to note that among the nine essential factors are factors which represent each of the three categories hypothesized at the outset of this research. That is  $X_4$  and  $X_{26}$  pertain to the distributed system network;  $X_8$ ,  $X_9$ ,  $X_{30}$  pertain to the application system; and  $X_{10}$ ,  $X_{11}$ ,  $X_{12}$ , and  $X_{14}$  pertain to the distribution policy.

The interpretation of coefficients describing the influence of qualitative variables is different than the interpretation of coefficients of quantitative variables. The coefficient of qualitative or indicator variables such as  $X_1$ ,  $X_2$ ,  $X_3$  and  $X_{10}$ ,  $X_{11}$ ,  $X_{12}$  describe the impact of change to that level from another level. The interpretation of coefficients modifying quantitative variables is traditional. That is, a positive coefficient corresponds to a direct relationship with the response variable and a negative coefficient designates an inverse relationship. It should be pointed out, however, in this research that the higher the value of the response variable the worse the performance. A strong inverse relationship between the variable and response is "good." Considerable caution should still be exercised in interpreting regression coefficients because regression does not imply causality. That is, there may be a strong correlative relationship between the factors which results in a significant regression, but the factors may not be related in a cause and effect fashion (20).

Table 19 continued. Coefficients for Variables in 10 Best Subsets Models

VAR X ( )	COEFFICIENTS									
	MODEL - 1	MODEL - 2	MODEL - 3	MODEL - 4	MODEL - 5	MODEL - 6	MODEL - 7	MODEL - 8	MODEL - 9	MODEL - 10
1	-	-	-	-	-	-	-	-	-	14.852
2	48.314	47.990	48.116	48.116	42.225	42.225	42.225	42.225	42.225	99.544
3	43.180	42.196	43.004	43.004	36.314	36.314	36.314	36.314	36.314	71.472
4	-29.091	-23.736	-28.638	-28.638	-31.943	-31.943	-31.943	-31.943	-31.943	-20.492
5	-	0.001	0.001	0.001	-	-	-	-	-	-
6	-19.687	-	-19.722	-19.722	-12.471	-12.471	-12.471	-12.471	-12.471	-20.001
7	7.471	-	7.422	7.422	3.154	3.154	3.154	3.154	3.154	7.366
8	1.214	-	1.219	1.219	0.777	0.777	0.777	0.777	0.777	1.229
9	-43.276	-	-43.324	-43.324	-50.688	-50.688	-50.688	-50.688	-50.688	-45.938
10	28.259	26.314	27.950	27.950	16.330	16.330	16.330	16.330	16.330	26.297
11	77.730	69.766	77.653	77.653	73.638	73.638	73.638	73.638	73.638	76.891
12	67.464	68.523	67.399	67.399	56.557	56.557	56.557	56.557	56.557	65.042
13	0.579	0.608	0.574	0.574	0.632	0.632	0.632	0.632	0.632	0.551
14	48.416	38.632	49.463	49.463	45.075	45.075	45.075	45.075	45.075	48.088
15	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
16	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
17	-	-0.006	-	-	-	-	-	-	-	-0.030
18	-	-	-	-	-	-	-	-	-	-
19	-456.723	-478.352	-454.700	-454.700	-464.700	-464.700	-464.700	-464.700	-464.700	-525.912
20	96.453	0.0	94.036	94.036	87.769	87.769	87.769	87.769	87.769	92.379
21	106.012	99.745	106.891	106.891	120.565	120.565	120.565	120.565	120.565	109.693
22	180.294	86.812	180.208	180.208	74.569	74.569	74.569	74.569	74.569	178.397
23	-	-	-	-	43.417	43.417	43.417	43.417	43.417	-
24	-8.343	-	-8.241	-8.241	-5.521	-5.521	-5.521	-5.521	-5.521	-8.332
25	10.540	13.336	10.297	10.297	9.091	9.091	9.091	9.091	9.091	10.704
26	-0.007	-0.012	-0.007	-0.007	-0.013	-0.013	-0.013	-0.013	-0.013	-0.007
27	-	0.044	-	-	0.053	0.053	0.053	0.053	0.053	-
28	-	-76.395	-	-	-154.000	-154.000	-154.000	-154.000	-154.000	-
29	-6.275	-10.657	-6.503	-6.503	-0.001	-0.001	-0.001	-0.001	-0.001	-5.167
30	-0.001	-0.001	-0.001	-0.001	-0.001	-0.001	-0.001	-0.001	-0.001	-0.001
31	-144.244	-83.409	-144.095	-144.095	-	-	-	-	-	-142.038
32	-	3.348	-	-	-	-	-	-	-	-

VARIABLE X { }	COEFFICIENTS				
	MODEL - 1	MODEL - 2	MODEL - 3	MODEL - 4	MODEL - 5
1	-	-	-	-29,706	-29,934
2	-	-	-	-	-
3	-	-	-	-	-
4	-20,793	-20,783	-22,407	-23,558	-24,001
5	-	0,001	0,001	0,001	-
6	-	-	-	-	-
7	-	-	-	-	-
8	1,149	1,149	0,982	0,965	0,861
9	-23,276	-23,146	-50,523	-53,975	-53,780
10	46,484	46,653	22,008	12,777	11,921
11	76,547	77,101	75,362	74,753	74,570
12	77,873	78,371	60,792	53,856	53,498
13	0,686	0,686	0,624	0,611	0,612
14	37,917	39,473	41,099	47,884	47,315
15	-	-	-	-	-
16	-	-	-	-	-
17	-	-	-	-	-
18	-	-	-	-	-
19	-282,503	-282,094	-307,898	-322,706	-323,400
20	-	-	-	-	-
21	-	-	-	-	-
22	49,342	49,993	102,631	116,506	125,148
23	-	-	-	-	-
24	-	-	55,479	68,253	69,140
25	12,343	12,340	8,540	7,367	7,520
26	-0,006	-0,006	-0,013	-0,015	-0,015
27	-	-	0,060	0,072	0,072
28	-94,383	-94,860	-192,765	-181,853	-179,748
29	-38,251	-37,380	-	-	-
30	-0,001	-0,001	-0,001	-0,001	-0,001
31	-	-	-	-	-
32	-	-	95,821	60,290	53,804

Since the units of measurement for the quantitative variables differ greatly, the magnitude of any given coefficient should not be construed solely as an indicator of influence. The units of measurement of any coefficient are the units of the response variable divided by the units of the regressor variable. The coefficients are determined jointly and serve to normalize variable values as well as measure the impact of individual variables on the response. That is, a coefficient indicates the influence of a given factor when that factor is considered simultaneously with all of the other factors in the model.

A regression coefficient measures the expected change in performance per unit change in the regressor, given that the levels of the other regressors in the model remain constant. This can obscure our understanding of the role of individual factors in the model if inferences are made about these factors in isolation. That some or many of the relationships may not transfer from the composite model to the single factor situation should be understood. Since the models described here comprise on the average 20 variables, it is of particular importance that this caution be observed. A minimum of 15 variables are required to explain performance when all the control factors are being manipulated.

#### Selection of Explanatory and Prediction Models

Table 20 presents a rank evaluation of the ten best subset models based on explanatory  $R^2$  and prediction  $R^2$ . With regard to quality of fit, all ten models are equivalent. It is apparent that the best explanatory models and the best prediction models do not coincide. The prediction  $R^2$  is in all instances lower than the explanatory  $R^2$ . This

is to be expected. An explanatory model is one which provides an adequate fit to the data on which it was built, and in which the regression coefficients are reasonable estimates of the effects of the predictor variables. A model that is a good predictor is generalizable; that is, it provides reasonable predictions of fresh data not used in the parameter estimation process.

The models developed in this research are all linear. They serve a factor screening function and as such perform very well. Obviously higher  $R^2$  values could be obtained if polynomial or other nonlinear models were fit. Such increases in model complexity are warranted only when they are grounded in physical reasons outside the data. This is certainly not the case here, as there is little, if any, underlying theory connecting the factors studied in this research to the performance response variables. Furthermore, in an experiment with only two levels of most factors such as this one, polynomial or nonlinear models are not meaningful.

The degree to which a model is satisfactory as a descriptor makes no implication concerning its generality. When a fitted model is applied to new data, it is unlikely to predict the fresh data as well as it fits the estimation data. Since the model is fit to the estimation set using least squares, it is, in some sense, an optimal fit for that data. Optimality here is unique to the estimation data set. Generally, a model which is 80 to 90 percent as satisfactory in prediction as it is in explanation is considered "acceptable" (20). Model 7 is only 6.4% as good in prediction as it is in description. Model 8, on the other hand, is 82.6% as good a predictor as it is at

explaining the data on which it was built.

The determining factors for model selection are model adequacy, generality and ease of use. Before model selection, then, the latter two criteria should be considered. Since those models which rank highest on explanatory  $R^2$  are not the same as those that rank highest on the predictive scale, we know that the "best" models are not the most general. A decision must in this case be made to either 1.) use two separate models for description and prediction or 2.) use a single model which compromises between these applications. If two models are to be used, the most likely choices would be the models which rank highest on the two adequacy scales. These would be Model 7 for description and Model 6 for prediction. If a single model is to be chosen the most likely candidates are Models 10 and 8. The difference between Models 10 and 8 on the explanatory scale is +.004. The difference between Models 10 and 8 on the prediction scale is -.05272. This difference in predictive capability suggests that Model 8 would be preferable to Model 10 as a general model. In fact, Model 8, as stated above, is 82.6 percent as good in prediction as description.

The remaining major consideration for model selection is ease of use. Our focus here will be limited to the four highest ranking models on the two  $R^2$  scales. The consecutive numbering of Models 1 through 10 corresponds to their ordering with regard to number of regressors. Model 1 has 15 regressors while Model 10 has 26. Since the top four models in terms of explanatory or prediction  $R^2$  have at least 24 regressors, model size is not a determining factor toward ease of use.

To aid in choosing between two models or a single model, i.e. 7

and 6; or 8 or 10, we refer to Table 20 to compare the variables that comprise each model. From this table we see that the only variable which is in Model 8 which is not in Models 6 or 7 is  $X_6$ . The only variables which are in Model 10 which are not in Models 6 or 7 are  $X_1$  and  $X_{18}$ .  $X_1$  is an indicator variable, therefore, its presence does not affect ease of use.  $X_6$  refers to node communication capacity. It is a direct measure which is trivially obtained.  $X_{18}$ , available communications capacity after initial assignment, is indirect and consequently more difficult to measure or estimate. This variable, which is present in Model 10, is the only one which differentiates the choices of Models 6 and 7; or Model 8 or 10 on the basis of ease of use.

If a choice is to be made between Models 8 and 10, Model 8 would be chosen on the basis of adequacy, generality and ease of use. Model 6 is 82.6 percent as good a predictor as Model 7 is at explanation which is exactly the same generality rating as Model 8. The difference between Models 6 and 8 on explanatory  $R^2$  is 0.0026 and between Models 7 and 8 on prediction  $R^2$  is 0.0078. Thus, it appears that Models 6 and 7 or Model 8 are essentially equivalent with respect to all evaluation criteria. Since it is usually considered preferable to use one model rather than two when all other attributes are constant, Model 8 is selected for use as the most satisfactory model for operational survivability and performance.

Table 20. Rank Ordering of 10 Best Subset Models

MODEL NO.	EXPLANATORY $R^2$	MODEL NO.	PREDICTION $R^2$	TRIM
7	.8652	6	.71536	
9	.8647	8	.71252	
10	.8641	2	.70286	
8	.8626	10	.66264	
6	.8622	1	.53577	
4	.8559	9	.50232	
5	.8551	3	.47850	
3	.8457	4	.33184	
2	.8326	5	.33173	
1	.8317	7	.05569	

### Discussion of Model Components

Before discussing in specific the inference of model components, two unique aspects of this research should preface. First, it should be noted that the experiment is conducted on highly stressed distributed systems. That is, the conditions imposed were exaggerated in order to test multiple aspects of influence. These severe conditions on processing resources, application system demands, etc. were such that little modifications would force the system to failure. Some treatment combinations leave the distributed system so highly packed that after loss of a small percent of the network resources it is extremely difficult, no matter what distribution policy is imposed, to recover. While highly stressing the distributed system allows us to determine the importance of certain factors, it sometimes requires special understanding of model components.

The second preliminary remark pertains to definition of the regressor variables. As was indicated earlier, it is hard to measure many of the attributes of distributed systems which are used in this research. However, in light of this difficulty and the large amount of controversy which surrounds measurement of software, performance, and distributed systems, the models developed here show profound stability (24). The variables as described serve the model building process very well and as will be shown function in a very comprehensive fashion.

Now let us examine the role of the quantitative variables in the 10 best subset models.  $X_4$ , number of nodes in the distributed system, which is in all models, has a negative coefficient. This is interpreted to mean that the more nodes there are in the distributed

system the more likely that performance will be satisfactory. As can be seen in Table 19, inverse relationships of this type exist between a number of the regressor variables and the response. These instances are discussed below.

Also examination shows that some of the regressor variables have positive coefficients which indicate a direct relationship with the magnitude of the response variable. Remembering that an increase in the value of the response means performance is moving toward failure, we interpret strong positive relationships as having a detrimental effect on performance and consequently on operational survivability.  $X_7$ , number of application system modules, suggests that performance will degenerate as the number of application system modules increases.  $X_{14}$  simply states that failure to initially assign the application system to the distributed system makes satisfactory performance difficult.

$X_{20}$  represents the ratio of memory requirements to useable memory capacity. The positive coefficient here says that as the memory requirements approach the total available memory, the likelihood of satisfactory performance decreases. A similar observation is made for  $X_{21}$  which represents the ratio of communications requirements to useable communications capacity.  $X_{22}$  designates application system connectivity. Its relationship to the response states that the higher the level of application system connectivity the poorer the prospects for satisfactory performance. Each of these relationships seem reasonable and confirm some of our intuitions about distributed systems.

$X_6$  indicates that the greater the capability of nodes to communicate with other nodes the more likely performance will be satisfactory. Since in this experiment the capacity of all links are held constant the communication capacity is determined strictly as a function of number of links.

Given this relationship between survivability and number of links, one would also expect distributed system connectivity to be an influential factor. Distributed system connectivity is represented by  $X_{19}$ . As expected, this factor is found in all models and in all cases demonstrates a strong inverse relationship to response.

To demonstrate how potentially misleading it is to interpret individual regression coefficients in a multiple regression framework, let us consider the case of  $X_1$ ,  $X_2$  and  $X_3$  which together represent the four distributed system topologies. These topologies are star, ring, network, and array and are represented by indicator variables  $X_1$ ,  $X_2$ ,  $X_3$  as discussed in Chapter IV. The coefficients for models four and five indicate that a change from the base topology, a star, to the ring topology will have an improving affect on performance. Models six through 10 further indicate that a change from the star to either the network or array would have a detrimental effect on performance. Figure 7, however, shows that average performance actually improves, although perhaps slightly, by a change from the star to any other topology. The model coefficients indicate the effect of topology given all the other factors in the model. For example, given that distributed system connectivity is represented by two quantitative variables,  $X_6$  and  $X_{19}$ , as discussed above, it might appear less

striking that the qualitative variable, distributed system topology, X4 which also represents this same feature, has a less profound inference than expected.

DISTRIBUTED SYSTEM TOPOLOGY DISTRIBUTION POLICY	DISTRIBUTED SYSTEM TOPOLOGY				AVERAGE
	STAR	RING	NETWORK	ARRAY	
RANDOM	3.29	1.75	1.77	3.31	2.53
UNIFORM	3.00	3.08	3.25	3.00	3.08
PACKED	3.74	3.76	3.95	2.89	3.59
OPTIMAL SPARE	3.81	4.00	3.25	4.00	3.77
AVERAGE	3.46	3.15	3.06	3.30	

Figure 7. Average Performance Given for Different Distributed System Topologies and Distribution Policies



When examining the effect of distribution policy, it is observed that distribution policy in all cases is important. A change from the random distribution to any other distribution effects a noticeable positive influence on the coefficient for the factor representing the new distribution approach. While Figure 7 bears this out, it also shows that with only two exceptions performance based on distribution policy is fairly uniform. And, performance based on the intersection of distributed system topology and distribution policy is even more homogeneous. That these factors are important and that on direct observation they seem indistinguishable appear contradictory. However, once again, it must be recognized that the importance of these factors comes from their role in the model when operating with numerous other factors.

$X_{30}$ ,  $X_{31}$ , and  $X_{32}$  represent interactions between other regressor variables.  $X_{30}$  signifies the interaction among a number of application system related attributes, namely; number of application system modules, module processing requirements, module memory requirements and module communication requirements.  $X_{31}$  and  $X_{32}$  signify the interaction between final dispersion, application system connectivity and distributed system connectivity. The sign attached to interaction variables is not as important as relative magnitude of the coefficients and the signs and magnitudes of the main effects of the variables in the interaction. That these features are important to performance seems reasonable and supports the initial postulate of this research concerning the believed complexity of adequate models of operational survivability.

Several factors included in some of the 10 best subset models have negligible effects. Interestingly, these factors  $X_5$ ,  $X_{15}$ ,  $X_{16}$ ,  $X_{17}$ ,  $X_{18}$ ,  $X_{21}$  (for definitions see Table 16) all represent direct measures such as node processing speed and global resource capacity, i.e., memory, processing, communication. In absolute or simple form these measures are not very meaningful, however, as has been shown,  $X_{24}$  and  $X_{30}$ , and as is shown,  $X_{20}$ ,  $X_{21}$ , and  $X_{25}$ , these direct measures when considered in conjunction with other attributes of the system can be extremely influential.

Another aspect of model inference which merits a word of caution is interpretation of the signs of regression coefficients. Frequently the signs of regression coefficients will coincide with prior expectation. Most often this occurs when 1.) all necessary regressor variables are in the model, 2.) the relationship between the variables and the response is strong and 3.) the regressors are orthogonal. Models with supersets and subsets of these constraints often demonstrate this status with coefficient signs which are counter intuitive. Such inconsistencies usually are minor if they pertain to the less important factors in the model.

One of the most common causes of "wrong" signs is multicollinearity. Multicollinearity refers to the existence of intercorrelation between the regressor variables. The eleven factors involved in the  $2^{K-P}_V$  Fractional Factorial experiment used in this research are orthogonal. However, a number of other candidate regressor variables were analyzed during the model building process. Many of these variables were derived from the original factors and

represent that factor in a somewhat specialized context. When all regressor variables used in the ten fitted models are analyzed simultaneously moderately strong multicollinearity is indicated. The most obvious solution is, of course, to simply remove the regressors involved in the multicollinearity. Removal of these variables, however, would destroy the predictive character of the model. Wrong signs in regression problems often occur for other reasons. For example, the violating factor may not be varied over a sufficiently wide range or necessary companion variables may be missing. The latter condition happens because a regression coefficient is a measurement of partial effect and does not stand alone. This type of wrong sign condition can sometimes be "corrected" by a redefinition of the variable.

It is not necessary that the signs of coefficients be in agreement with prior expectation. The degree to which the factors fall short of fulfilling their combined role of explaining or predicting the response is reflected in the model adequacy evaluations. Adjustments which influence the direction of signs such that they concur with expectation may result in models with higher adequacy ratings. This does not imply, however, that sign concurrence will assure an increase in model adequacy or that a model that fits a set of estimation data will have intuitive appeal or be a useful predictor of new observations. For further discussion see Montgomery and Peck (20).

Inferences of some model components require thoughtful interpretation. Factors  $X_{20}$  and  $X_{28}$ , for example, are to some extent related, however, not enough to be determined redundant. Redundancy

would indicate that one of the variables could be removed without affecting the model. Here, we find that either both  $X_{26}$ , available processing capacity at the end of the subcase, and  $X_{28}$ , dispersion at the end of the subcase, are present in the model or just  $X_{26}$  is.  $X_{28}$  represents the final number of nodes over which the application system is distributed divided by the number of application system modules. The closer dispersion comes to being total, or one, the more likely performance is to be satisfactory. The potential for dispersion is, of course, related to the resources available. Thus, a somewhat collinear relationship between  $X_{26}$  and  $X_{28}$  is to be expected.

For purposes of inference an interrelationship exists between all the regressor variables that relate either directly or indirectly to dispersion.

$X_{23}$  represents dispersion after initial assignment. It is implied that the greater initial dispersion that exists the less likely performance will be satisfactory. The apparent contradiction between the relationship of  $X_{23}$  to the response and that of  $X_{28}$  to the response requires further investigation. The coefficient on  $X_{23}$  infers initial dispersion is "bad" and the coefficient on  $X_{28}$  infers final dispersion is "good." The phenomenon observed here results from the highly stressed nature of this experiment. That is, the design of the distributed systems tested was such that if a small number of nodes were lost, features of the system other than simply excess processing and memory capacity were required to make it survive. Thus, if the application system was initially dispersed there would be an increased probability that losing nodes would make it impossible to recover. If,

on the other hand, the application system was concentrated on a few nodes, the likelihood of losing valuable nodes would decrease resulting in a higher probability of survival.

This situation is further exhibited by  $X_{29}$ , criticality of the lost nodes. Criticality here is determined as the ratio of the sum of the connectivity of application modules on the lost nodes to the application system connectivity. The coefficient states that the closer the criticality ratio comes to one the more satisfactory performance. Although this relationship is worthy of further study, the following is offered as a possible explanation. Previously it was postulated that final dispersion has a constructive relationship to performance. This factor, however, seems to imply that when examining lost nodes concentration is desirable. It is possible that both of these conditions hold, however, one pertains to performance and the other to likelihood of satisfactory reconfiguration. The models infer that initial concentration of the application system is desirable. It follows, then, that concentration of that which is lost will facilitate recovery.

$X_{24}$  represents memory consistency, that is the number of application system modules divided by the number of application system modules that will "fit" on a node memory-wise. This says that as this ratio increases chances for survival improve. In other words, the closer the system can come to placing all the application system modules on a single node, the less likely that performance will be satisfactory. Once again, if we relate dispersion to performance and concentration to recovery, this inference is reasonable. The fewer the

application modules that will fit on a node the more likely that the application system will be dispersed.

$X_{25}$ , processor consistency, represents the ratio of number of application system modules to average number of modules which will "fit" on a node processor-wise. The positive coefficient here indicates that as this ratio increases performance degrades. Such a proposal is intuitive. The probability of satisfactory application system performance will increase directly with the ability to assign all processing to a single node. However, when capability falls short of that, the importance of distribution policy and connectivity may become dominant. It may not be that high processor consistency is detrimental when the ratio is greater than one but that in complex systems reconfiguration is difficult.

It is apparent that dispersion and concentration are companion concepts. Further analysis using some of these factors in a designed experiment so that their main effects and interactions may be more precisely estimated is desirable. Exercising these factors in less highly stressed experiments may be necessary. Experiments of this type should be useful in clarifying the relationship between these variables and the response.

The research documented in this dissertation demonstrates both the capability and significance of empirical investigation in distributed processing. The experimental results presented do not support conclusions drawn from prior analytical models (19). Merwin and Mirhakak's survivability index, for example, indicates that the most survivable distributed network is a star. That determination is

based on number of links to be traversed between any two network nodes. This research clearly shows that other factors such as potential for alternate routes, characteristics of the application system and distribution method are also strongly influential. When averaging responses based on topology and distribution policy all three other topologies tested fared better than the star topology. Differences are further highlighted when all model components are considered. These results bring into question the present capability of analytical models to represent complex problems of inexact sciences. It is also apparent, however, that analytical modeling may be appropriate for examination of specific model components such as those which can be expressed in totally quantitative terms. The three consistency measures fall into that category. Empirical methods with which to test analytical models are available. Used together, these methods should lead to a strongly quantitative understanding of survivability which can be used in the design of distributed systems and validated in field tests.

## CHAPTER VIII

### CONCLUSIONS AND RECOMMENDATIONS

One objective of this research was to enhance our understanding of operational survivability and performance and to make that understanding quantitative. The approach taken was an experimental one which used factor screening to give indication of variable importance. The objective was to develop models which are explanatory and would provide a foundation for future refinements rather than prescriptive. The second objective of this research was to demonstrate the applicability of traditional experimental design and regression analysis techniques to the field of computer science. The experiment and results documented in this dissertation support these objectives.

A factor screening experiment was conducted to determine whether any of a large set of candidate regressor variables were important to operational survivability and performance. Results demonstrate that a number of variables are, indeed, very influential and analysis shows their approximate level of importance. A two level factor screening design with a large number of variables was used in this research. Given this experimental approach, it is relatively unusual and encouraging that the design provides sufficient information on which to build ten linear explanatory models with  $R^2$  values in the range of .8. Further, the capability of several of these models to serve exceptionally well in a predictive role suggests that they provide a good foundation on which to build future refinements.

in influence. Each of these are directly measurable entities such as global memory, processing and communications capacity. The more important factors tended to be more complex or more indirectly derived. Examples are distributed system connectivity, application system connectivity and memory consistency. This finding further supports the initial proposal that operational survivability cannot be trivially indexed.

In summary, 32 candidate regressors are used in identifying the 10 best subset models. The coefficients of these regressors are approximately equivalent in sign and magnitude across models. All variables remain proportional with the introduction and removal of other variables, thereby demonstrating extreme stability. The explanatory adequacy of models built using these variables is in all instances in excess of .8 which is very acceptable for a factor screening experiment. The adequacy in prediction of these models ranges between -.29 and +.71 with some models predicting very well and others predicting very poorly. By constructing satisfactory explanatory and predictive models, this research demonstrates that the concept of operational survivability and performance as proposed can be expressed quantitatively. Further, it is shown that major factors include the distributed system network, application system and distribution policy as initially proposed.

In review we find that there are nine factors found in all models. These are number of nodes in the distributed system, distributed system connectivity, module memory requirements, module to module interaction frequency, distribution policy, percent nodes lost,

The models developed here through standard regression techniques make a number of statements about measurement of operational survivability and performance. The first and most important statement is that these attributes can, in fact, be described in a quantitative fashion. Next, they imply that certain factors are more important than others in determining the level of the response. Some of the most influential factors are distributed system connectivity, number of nodes, available processing capacity, distribution policy, application system connectivity and module memory requirements. The number of regressor variables required to achieve explanatory model adequacy levels of .8 is large. Large is here defined as between 15 and 26. The nine core variables found in all ten models have the expected sign and an obvious interpretation. The few instances in which signs do not concur with expectation occur in connection with peripheral or less important factors. These instances are well within acceptable bounds for research of the type conducted here. It is further shown that among the nine essential factors are factors which represent the three general categories hypothesized at the outset of this research. Also, it is demonstrated that no single category or pair of categories will adequately explain or predict operational survivability or performance. The three categories describe attributes of the distributed system network, application system and distribution policy.

Analysis of the experiment results supports the hypothesis that the factors necessary to adequately describe operational survivability would be large in number and non-trivial in observation. The ten best subset models included a number of factors which were nominal

initial assignment results, available processing capacity at the end of the subcase and the interaction of all application related variables.

Other factors which prove to be important and function in the models in an expected manner are number of application modules, node communication capacity, memory requirements ratio, communication requirements ratio, and application system connectivity. Some factors operating as expected given the highly stressed nature of the experiment conducted are initial and final dispersion; memory and processor consistency; and criticality of lost nodes.

Factors having negligible effect include node processing speed; global memory capacity; available processing, memory and communications capacity after initial assignment; and available communications capacity at the end of the subcase.

A number of propositions can be inferred from the analyses of Chapter VII. Some of these are not unexpected. Others, however, are somewhat surprising, and we offer them as hypotheses which can be further explained experimentally. While plausible explanation can be offered to support each of these hypotheses, there are also apparently plausible settings in which the hypotheses may fail. Both confirming instances and refutations of the hypotheses point the way toward further experimentation. That is, for each of the 10 hypotheses listed below we present a possible mechanism to explain the effect which is apparently being observed. We then give a brief indication of situations in which the hypothesis may fail; the appropriate experimental setting for dealing with the hypothesis should lie within these limits.

1.) The more nodes there are in the distributed network, the more likely that performance is satisfactory. It seems very likely that the distributed systems in which we are most interested satisfy this property, that is the more nodes there are in a distributed network configuration the more likely there will be slack or excess resource capacity which can be used if other resources are lost. However, given a ring network configuration with communication links traveling in only one direction, the loss of a single node will destroy the network no matter how many nodes it contains. Likewise, this is true for a star configuration if the central node is the node lost.

2.) As the memory requirements approach the total available memory, the likelihood of satisfactory performance decreases. Given an application system which is distributed over the nodes of a distributed network, it is reasonable to conclude that as the demands on memory approach the memory limit of the network the more likely additional resource losses will have a detrimental effect on survivability due to constraints on reconfiguration options. On the other hand, it is apparent that as the distributed network decreases so too does the available memory until finally the memory available is only that on a single node. Further, it is possible that the memory requirements of the application system are extremely low and fit well within the memory capacity of a single node, however, the processing demands

exceed the capabilities of the processor. In this case the memory requirements to availability ratio has no relationship to survivability.

3.) As the module interaction or communications requirements approach the total available communications capacity, the likelihood of satisfactory performance decreases. It is not difficult to envision a number of network configurations in which the options for satisfactory reassignment of application modules decreases as the communications demands of the application system approach the communication limit of the distributed system. However, if the interaction requirements of application modules is such that those modules having the highest interaction can always be placed on a single node this relationship may not hold. Also, if the network configuration is such that two large subnetworks are connected by a bridge and the application system is split such that a large portion of the module to module interactions must traverse the bridge, the performance may decrease even though the available communication capacity is high.

4.) The higher the distributed network connectivity, the greater its probability of survival. Research in network survivability and routing support our basic intuition that in general the larger the number of alternate routes available for nodes to communicate with other nodes the greater the likelihood that an application system spread

over several nodes will be able to continue to adapt to increased node losses. Special cases can be identified to which this general statement does not apply. For example, if a network comprising nodes and links of low capacity or nodes and links which are nearly saturated is highly connected and the distribution/redistribution policy is such that tasks are dynamically reassigned to "optimize" node and link utilization, the fact that the options are numerous may be a drawback. In other instances high network connectivity may be irrelevant to survivability. For example, if a network is highly connected but the application system to be executed on it comprises only two modules, the degree of network connectivity may be of negligible importance.

5.) Failure to properly assign the application system to the distributed network initially makes satisfactory or degraded performance difficult. The complexity of mapping an application topology onto a network topology increases with the size and connectivity of the two graphs to be mapped. Thus, when an application system is assigned to the distributed system in such a way that it does not meet performance requirements, adjustments to correct the problem require additional sophistication on the part of the redistribution strategy. That is, once the problem of unsatisfactory performance is detected, the cause must be determined and a solution found. Depending on the distribution/redistribution policy the solution space for correction is often

more constrained than the initial solution space. On the other hand, if the distribution/redistribution algorithm is an adaptive one that examines different distribution options to determine their effect, then unsatisfactory initial allocations may be more useful or informative than satisfactory distributions. Unsatisfactory distributions may provide insight into worst case conditions.

6.) Performance degenerates as the number of application modules increases. We are essentially postulating that as the number of application modules increases the task of assigning and reassigning them in such a way that performance is satisfactory becomes increasingly more complicated. This is particularly true if the modules have high interaction requirements and few options for assignment due to module size or network configuration constraints. To see how such a mechanism could fail to hold, let the application system be of size  $N$ . The choice exists to either have five modules of size  $N/5$  or 20 modules of size  $N/20$  on a 10 node network, any node of which can accommodate an  $N/4$  size module. It is clear that having more modules offers more flexibility and possibly more opportunity for satisfactory assignment. Here, the larger number of small modules potentially fit on four nodes and could disperse to 10 nodes. The larger modules need at minimum five nodes. Maximum dispersion for the larger modules is also five nodes.

7.) The higher the level of application system connectivity the poorer the prospects for satisfactory performance. Here again, there is indication that high software system complexity will influence survivability. Given an application system for which module requirements nearly correspond to individual node capabilities, high connectivity may require a one for one mapping of the application system onto the distributed network. If such a mapping can be constructed initially it is unlikely that it can be maintained with increasing node losses. There are also instances in which software complexity may have little or no effect on survivability and performance. For example, an application system can be highly connected but have module to module interaction frequencies so low that as long as there is a path from any module to any other module the interaction demands can be met.

8.) The greater initial dispersion the less likely performance will be satisfactory. Given a distributed network of high or low connectivity it is not difficult to find situations in which the greater initial dispersion the more difficult recovery due to reduced reconfiguration options. Depending on the distribution/redistribution approach, however, it may be that the greater initial dispersion the fewer application modules to be reassigned after the loss of any single node. This would indicate that initial dispersion has a positive influence on performance.



9.) The greater final dispersion the more likely performance will be satisfactory. Corresponding to the previous inference, the greater the initial flexibility in the system the greater the opportunities for subsequent reconfiguration. The greater final dispersion the less saturated the system resources. While this argument may hold it is also possible to imagine application systems for which the postulate may not be true. For example, the greater final dispersion the less likely that highly connected application systems with high average module to module interaction frequency will be assigned such that their performance requirements can be met.

10.) The larger the proportion of highly connected application modules on the nodes lost the greater the likelihood of survival. Like hypothesis (8), this hypothesis concerns the effects of possible reconfiguration options. The effect in this case is simply one of removing logical dependencies: as dependencies are removed, the remaining nodes become (if they still meet the application requirements) autonomous and this can be exploited in assigning the remaining resources. Again, special cases can be described for which this argument does not hold. One such case is that in which the modules to be reassigned are highly connected but the network onto which they are to be placed is heavily saturated and available resources are widely dispersed. Given these conditions it is likely that performance will

degrade rather than improve. Thus, it appears having the flexibility to reassign all modules having the most severe interaction constraints can facilitate successful reassignment given the network resources are not heavily saturated.

The experimental and modeling techniques used in this dissection represent an initial step in developing a measurement instrument for operational survivability in gracefully degrading distributed processing systems. Further refinements in this measurement tool should facilitate more precision in its explanatory and predictive capability. Other recommendations for future research fall into two categories. These categories are 1.) more extensive use of the simulator as an experimental device and additions to its current capabilities and 2.) experimentation to clarify the operation of specific factors. SURSIM is a fairly general purpose simulator. The parameter levels used for this research designate selections made for this factor screening experiment. They do not represent limitations of the simulator. Using the variables under its control the simulator can generate a virtually unlimited number of treatment combinations. This provides the capability to focus future experimentation on some single or small set of factors while fixing the context environment with appropriate constants. There are features of the simulator which were not exercised in this experiment. One of these was to vary the capacities of the communication links. Also, heterogeneous distributed systems can be described and accommodated by the simulator manipulation and evaluation routines. Among the possible additions to the simulator

are the capability to represent multiple communications links between nodes; limitation on the availability of software; node and link vulnerability and criticality; and a larger number of distribution policies.

Future research which is likely to be productive includes experimentation on factors related to dispersion, connectivity and distribution policy. Designed experiments which focus on the control of these factors should improve our understanding of their direct and indirect operation. The introduction of more interaction variables could also be helpful. Also, exercising the factors examined in less highly stressed experiments may be meaningful.

The research conducted here identifies the variables important to operational survivability and to some extent tells how large changes in these important variables affect the response. Future experimentation which provides either a large number of factor levels or finer granularity in possible variable values should permit greater resolution in the simulator results and their subsequent application. The results presented in this dissertation demonstrate the applicability of traditional experimentation and regression analysis in the field of computer science as well as the feasibility of measurements which can serve as measurements for distributed systems. The models developed represent a promising initial step in the quantification of operational survivability as it applies to gracefully degrading distributed processing systems.

## APPENDICES

## APPENDIX A

DESCRIPTION OF DATA USED IN  
DESIGNED EXPERIMENTS

A description of data used in the 128 designed experiments is presented below.

1. Factor  $Z_1$  - Distributed System Topology:

Four different topologies are used in this experiment. They are a star, ring, network and array. Examples of these topologies for four and 10 node networks are presented in Figures A-1 through A-4.

2. Factor  $Z_2$  - Number of Nodes:

Two different size distributed networks are used in this experiment. These comprise 4 and 10 nodes respectively.

3. Factor  $Z_3$  - Node Processing Speed:

Two node processing speeds are used in this experiment. They are 500 kilo operations per second or 500 kps and 10 million operations per second or 10 mops.

4. Factor  $Z_4$  - Node Memory Capacity:

Two different node memory capacities are used. These are 128 kilobytes or 128 kbytes and 2 megabytes or 2 mbytes.

5. Factor  $Z_5$  - Connectivity of Applications System:

Application systems with high and low connectivity are used. The topology of these systems for the 4 and 16 node application systems used in this experiment are presented in Figures A-5 through A-8.

6. Factor  $Z_6$  - Number of Application Modules:

Two different quantities of application system modules are used in this experiment. They are 4 and 16 modules respectively.

7. Factor  $Z_7$  - Average Module Processing Requirements:

Application module processing requirements are computed as .1 or .5 of the node processing capacity after total network processing capacity is divided by the quantity of application system modules.

8. Factor  $Z_8$  - Average Module Memory Requirements:

Application module memory requirements are computed as .1 or .8 of the node memory capacity after total network memory capacity is divided by the quantity of application system

modules.

9. Factor  $Z_9$  - Average Module to Module Interaction Frequency:

Two levels of interaction frequency are used. These are high interaction frequency, which is computed as 50% of the average module processing requirements; and low interaction frequency, which is computed as 1% of the average module processing requirements. These frequencies are expressed in thousands of messages or packets sent per execution of an application module.

10. Factor  $Z_{10}$  - Distribution/Redistribution Policy:

Application system modules are assigned to the distributed system topologies according to one of four possible graph mapping algorithms. The algorithms used in this experiment are defined as follows.

Random Distribution - Application system modules are randomly assigned to processors. If the application module and communication burden will not fit at the node selected another random assignment will not be made. This will be repeated until all modules have been assigned to node. Should this approach fail to construct a map, the simulator in its present form will not attempt to degrade or reconfigure the system.

Uniform Distribution - Application system modules are assigned to nodes such that each node has as near the same operating demands as possible. This type of distribution is relatively easy to implement in central processor or master/slave type systems. Distributed systems in which global information about the system is available to each node must take into account the overhead burden this will place on the system resources. The overhead burden is dependent upon the size of the distributed system and timeliness of information required, i.e., frequency of update. (In distributed systems with high capability nodes, the impact of this update activity may be negligible. For distributed systems with a large number of low capability nodes, this burden is possibly very significant.) For the simulation under discussion such overhead burden will not be a factor; however, given some rule to be used to determine overhead burden incorporation into the model would be possible.

Packed Distribution - Application system modules are assigned to a designated processor until it reaches maximum capacity after which point modules are assigned to the next (nearest) processor, etc. If multiple processors are one communication link away the next node to be packed will be randomly chosen.

Optimal Spare Distribution - Application system modules are

assigned to the distributed processing system in such a way that each node being assigned application tasks has a spare queue indicating the sequence of backup or spare nodes which will be activated should the former fail. If insufficient nodes are available to provide every node with a spare, spares will be given to the nodes with application modules having the highest criticality ranking. Other "spares" may be shared by nodes executing lower criticality software. The concept of optimal-spare will become more complex and perhaps yet more meaningful when the vulnerability attribute is incorporated into the model.

11. Factor  $Z_{11}$  - Percent Nodes Eliminated:

Four different ranges of percent node elimination are used. These are

- 1 - 10%
- 11 - 30%
- 31 - 50%
- 51 - 80%

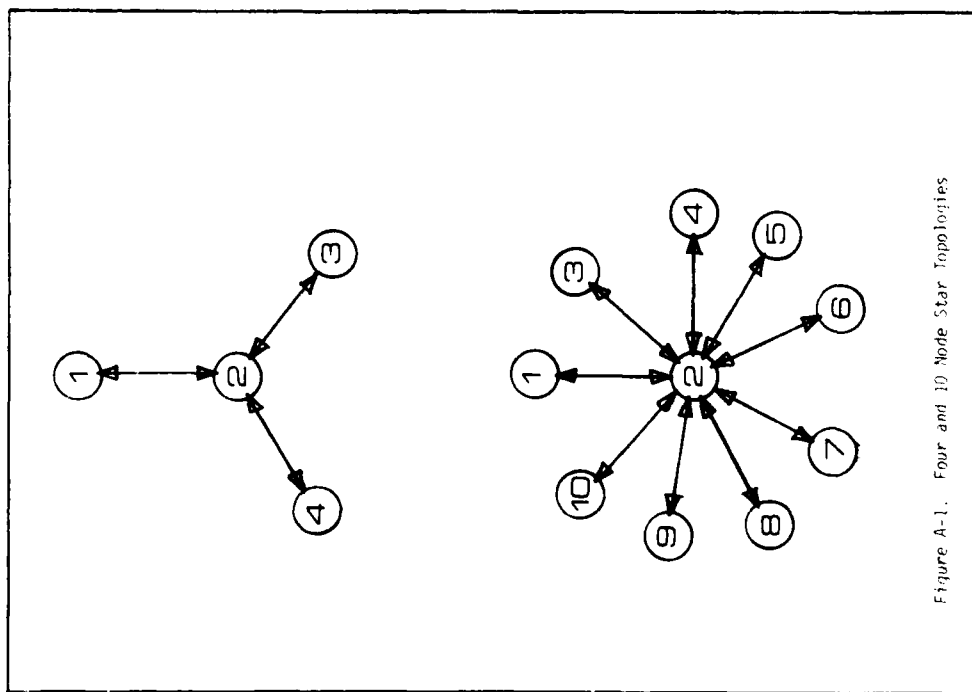


Figure A-1. Four and 10 Node Star Topologies

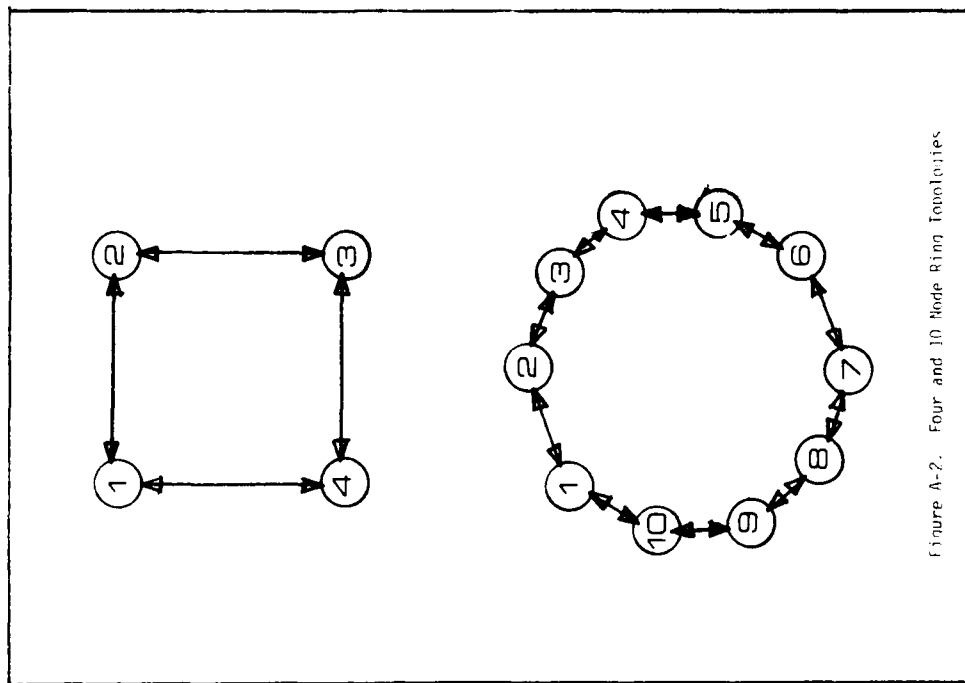
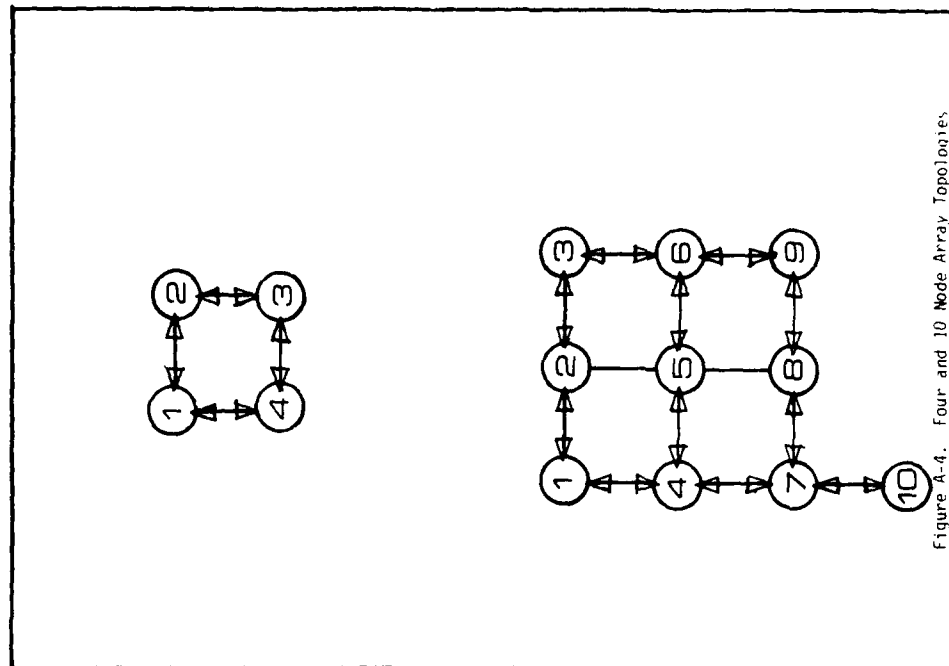
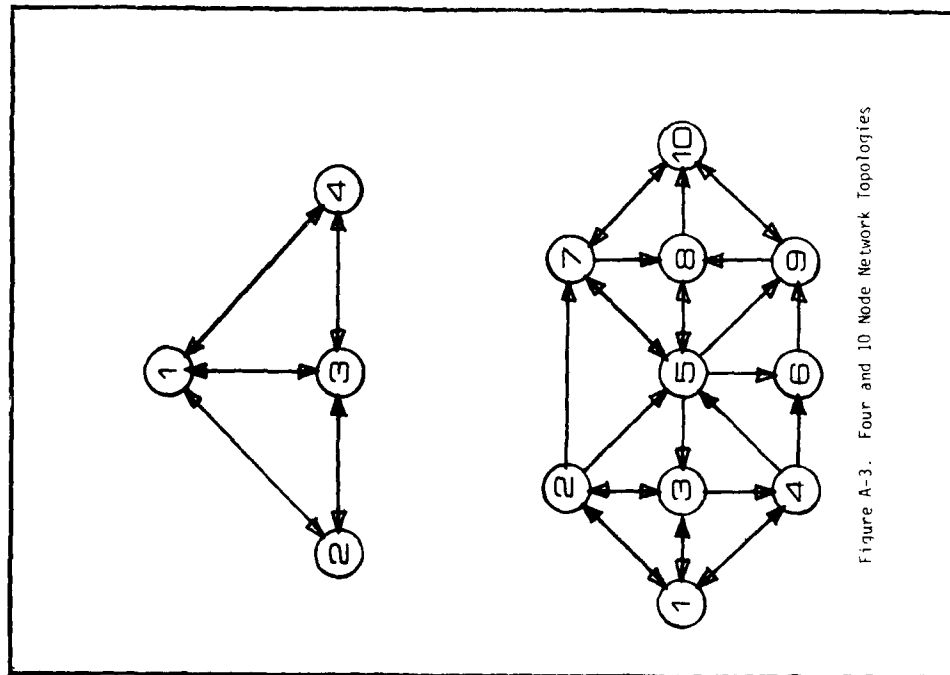


Figure A-2. Four and 10 Node Ring Topologies



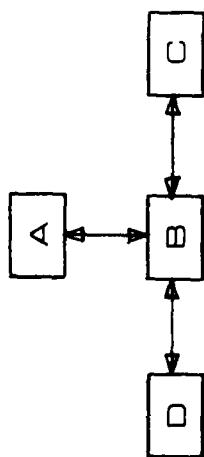


Figure A-5. Four Module Application System - Low Connectivity

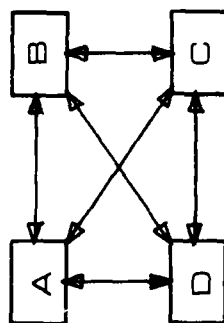


Figure A-6. Four Module Application System - High Connectivity

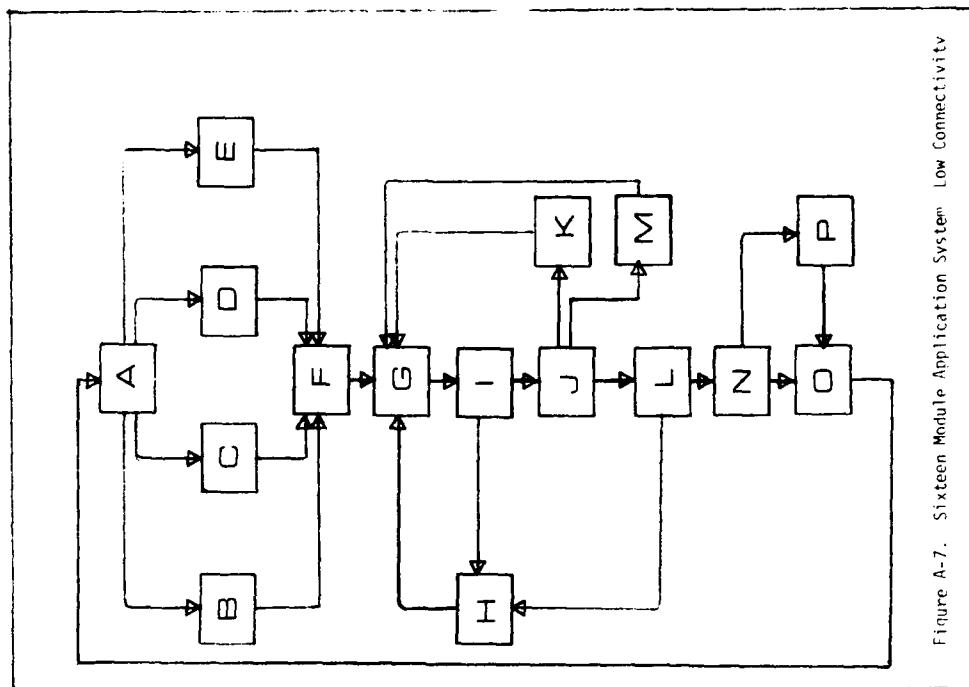


Figure A-7. Sixteen Module Application System - Low Connectivity

APPENDIX B  
TEN OPTIMAL SUBSET MODELS

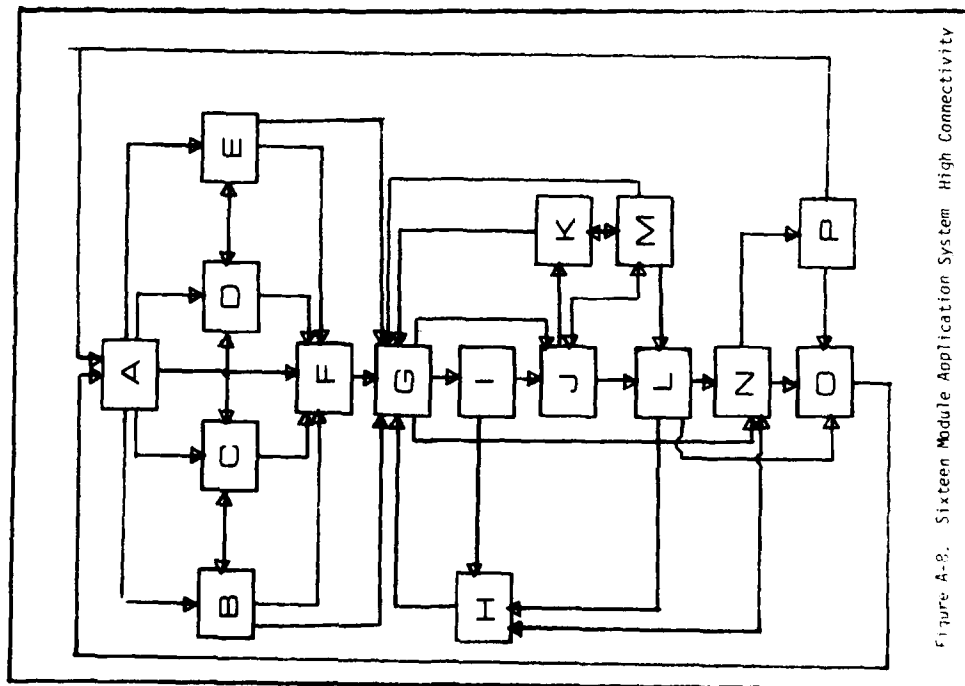


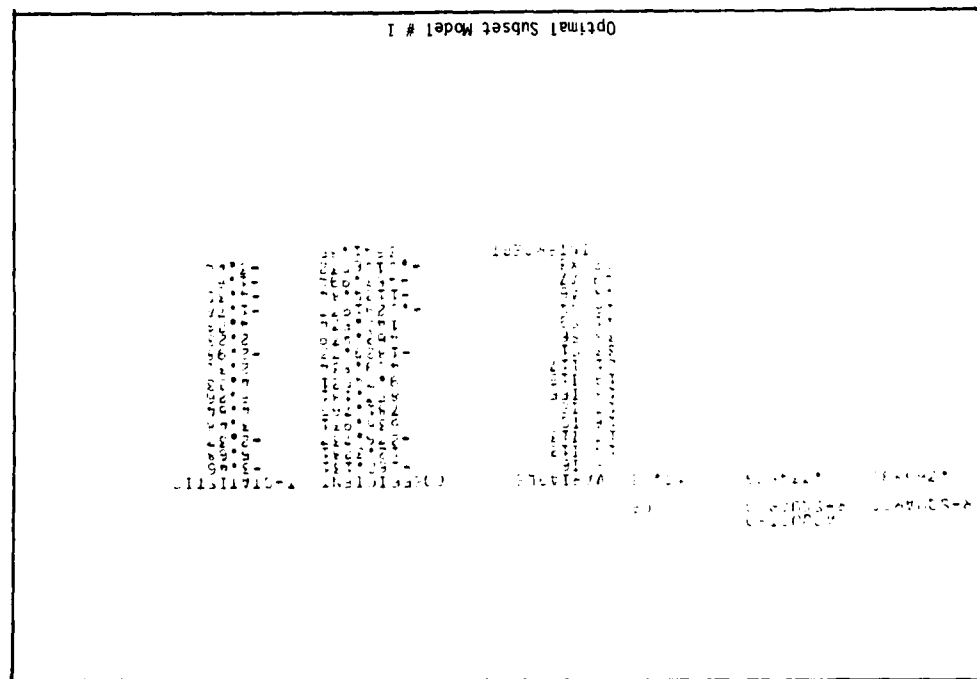
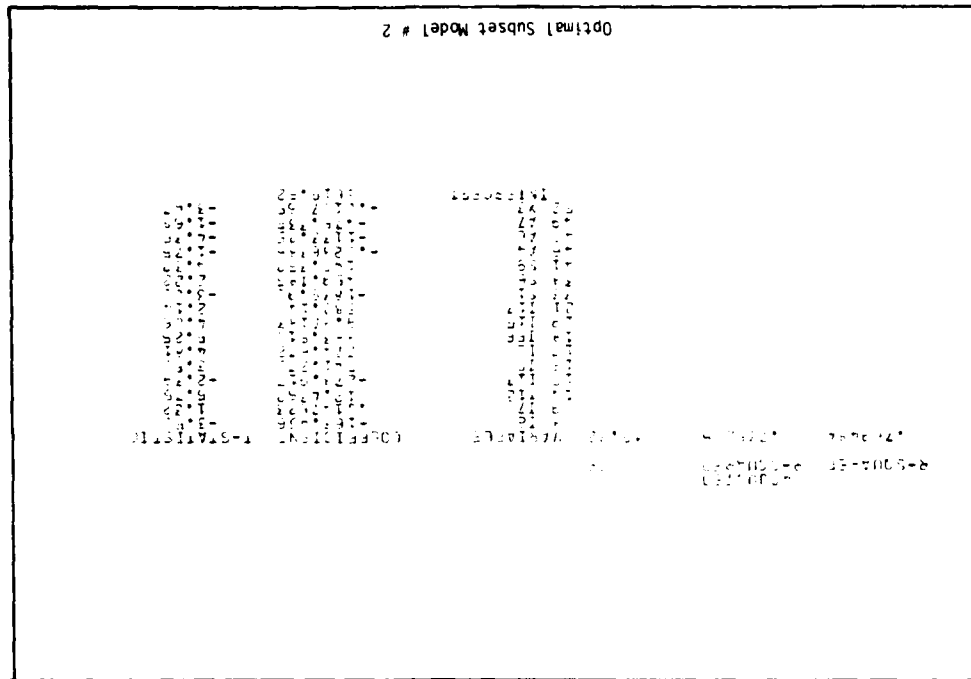
Figure A-8. Sixteen Module Application System High Connectivity



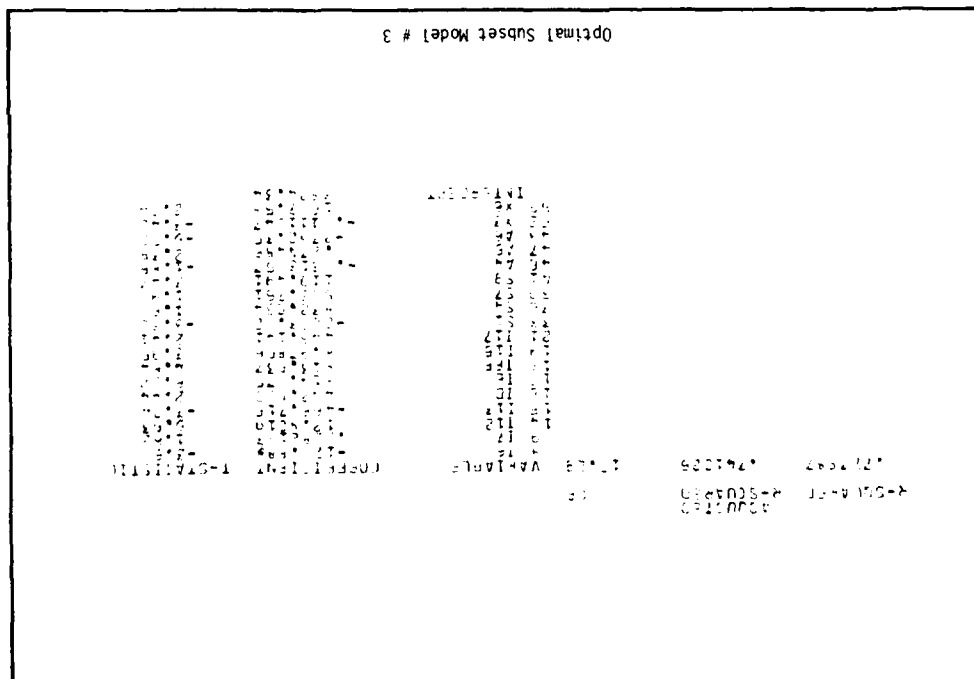
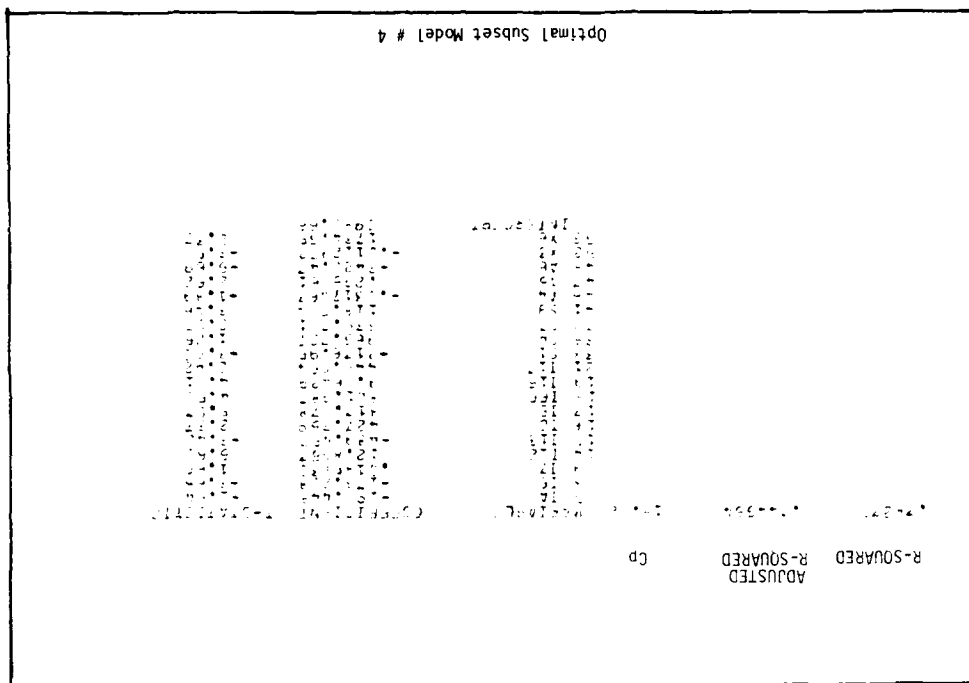
## Variable Key

Factor  
Description

11	12	13	14	15	16	17	18	19	20	21	22	23	24
115	117	119	121	123	125	127	129	131	133	135	137	139	141
143	145	147	149	151	153	155	157	159	161	163	165	167	169
171	173	175	177	179	181	183	185	187	189	191	193	195	197
199	201	203	205	207	209	211	213	215	217	219	221	223	225
227	229	231	233	235	237	239	241	243	245	247	249	251	253
255	257	259	261	263	265	267	269	271	273	275	277	279	281
283	285	287	289	291	293	295	297	299	301	303	305	307	309
311	313	315	317	319	321	323	325	327	329	331	333	335	337
339	341	343	345	347	349	351	353	355	357	359	361	363	365
367	369	371	373	375	377	379	381	383	385	387	389	391	393
395	397	399	401	403	405	407	409	411	413	415	417	419	421
423	425	427	429	431	433	435	437	439	441	443	445	447	449
451	453	455	457	459	461	463	465	467	469	471	473	475	477
479	481	483	485	487	489	491	493	495	497	499	501	503	505
507	509	511	513	515	517	519	521	523	525	527	529	531	533
535	537	539	541	543	545	547	549	551	553	555	557	559	561
563	565	567	569	571	573	575	577	579	581	583	585	587	589
591	593	595	597	599	601	603	605	607	609	611	613	615	617
619	621	623	625	627	629	631	633	635	637	639	641	643	645
647	649	651	653	655	657	659	661	663	665	667	669	671	673
675	677	679	681	683	685	687	689	691	693	695	697	699	701
703	705	707	709	711	713	715	717	719	721	723	725	727	729
731	733	735	737	739	741	743	745	747	749	751	753	755	757
759	761	763	765	767	769	771	773	775	777	779	781	783	785
787	789	791	793	795	797	799	801	803	805	807	809	811	813
815	817	819	821	823	825	827	829	831	833	835	837	839	841
843	845	847	849	851	853	855	857	859	861	863	865	867	869
871	873	875	877	879	881	883	885	887	889	891	893	895	897
899	901	903	905	907	909	911	913	915	917	919	921	923	925
927	929	931	933	935	937	939	941	943	945	947	949	951	953
955	957	959	961	963	965	967	969	971	973	975	977	979	981
983	985	987	989	991	993	995	997	999	1001	1003	1005	1007	1009
1011	1013	1015	1017	1019	1021	1023	1025	1027	1029	1031	1033	1035	1037
1039	1041	1043	1045	1047	1049	1051	1053	1055	1057	1059	1061	1063	1065
1067	1069	1071	1073	1075	1077	1079	1081	1083	1085	1087	1089	1091	1093
1095	1097	1099	1101	1103	1105	1107	1109	1111	1113	1115	1117	1119	1121
1123	1125	1127	1129	1131	1133	1135	1137	1139	1141	1143	1145	1147	1149
1151	1153	1155	1157	1159	1161	1163	1165	1167	1169	1171	1173	1175	1177
1179	1181	1183	1185	1187	1189	1191	1193	1195	1197	1199	1201	1203	1205
1207	1209	1211	1213	1215	1217	1219	1221	1223	1225	1227	1229	1231	1233
1235	1237	1239	1241	1243	1245	1247	1249	1251	1253	1255	1257	1259	1261
1263	1265	1267	1269	1271	1273	1275	1277	1279	1281	1283	1285	1287	1289
1291	1293	1295	1297	1299	1301	1303	1305	1307	1309	1311	1313	1315	1317
1319	1321	1323	1325	1327	1329	1331	1333	1335	1337	1339	1341	1343	1345
1347	1349	1351	1353	1355	1357	1359	1361	1363	1365	1367	1369	1371	1373
1375	1377	1379	1381	1383	1385	1387	1389	1391	1393	1395	1397	1399	1401
1403	1405	1407	1409	1411	1413	1415	1417	1419	1421	1423	1425	1427	1429
1431	1433	1435	1437	1439	1441	1443	1445	1447	1449	1451	1453	1455	1457
1459	1461	1463	1465	1467	1469	1471	1473	1475	1477	1479	1481	1483	1485
1487	1489	1491	1493	1495	1497	1499	1501	1503	1505	1507	1509	1511	1513
1515	1517	1519	1521	1523	1525	1527	1529	1531	1533	1535	1537	1539	1541
1543	1545	1547	1549	1551	1553	1555	1557	1559	1561	1563	1565	1567	1569
1571	1573	1575	1577	1579	1581	1583	1585	1587	1589	1591	1593	1595	1597
1599	1601	1603	1605	1607	1609	1611	1613	1615	1617	1619	1621	1623	1625
1627	1629	1631	1633	1635	1637	1639	1641	1643	1645	1647	1649	1651	1653
1655	1657	1659	1661	1663	1665	1667	1669	1671	1673	1675	1677	1679	1681
1683	1685	1687	1689	1691	1693	1695	1697	1699	1701	1703	1705	1707	1709
1711	1713	1715	1717	1719	1721	1723	1725	1727	1729	1731	1733	1735	1737
1739	1741	1743	1745	1747	1749	1751	1753	1755	1757	1759	1761	1763	1765
1767	1769	1771	1773	1775	1777	1779	1781	1783	1785	1787	1789	1791	1793
1795	1797	1799	1801	1803	1805	1807	1809	1811	1813	1815	1817	1819	1821
1823	1825	1827	1829	1831	1833	1835	1837	1839	1841	1843	1845	1847	1849
1851	1853	1855	1857	1859	1861	1863	1865	1867	1869	1871	1873	1875	1877
1879	1881	1883	1885	1887	1889	1891	1893	1895	1897	1899	1901	1903	1905
1907	1909	1911	1913	1915	1917	1919	1921	1923	1925	1927	1929	1931	1933
1935	1937	1939	1941	1943	1945	1947	1949	1951	1953	1955	1957	1959	1961
1963	1965	1967	1969	1971	1973	1975	1977	1979	1981	1983	1985	1987	1989
1991	1993	1995	1997	1999	2001	2003	2005	2007	2009	2011	2013	2015	2017
2019	2021	2023	2025	2027	2029	2031	2033	2035	2037	2039	2041	2043	2045
2047	2049	2051	2053	2055	2057	2059	2061	2063	2065	2067	2069	2071	2073
2075	2077	2079	2081	2083	2085	2087	2089	2091	2093	2095	2097	2099	2101
2103	2105	2107	2109	2111	2113	2115	2117	2119	2121	2123	2125	2127	2129
2131	2133	2135	2137	2139	2141	2143	2145	2147	2149	2151	2153	2155	2157
2159	2161	2163	2165	2167	2169	2171	2173	2175	2177	2179	2181	2183	2185
2187	2189	2191	2193	2195	2197	2199	2201	2203	2205	2207	2209	2211	2213
2215	2217	2219	2221	2223	2225	2227	2229	2231	2233	2235	2237	2239	2241
2243	2245	2247	2249	2251	2253	2255	2257	2259	2261	2263	2265	2267	2269
2271	2273	2275	2277	2279	2281	2283	2285	2287	2289	2291	2293	2295	2297
2299	2301	2303	2305	2307	2309	2311	2313	2315	2317	2319	2321	2323	2325
2327	2329	2331	2333	2335	2337	2339	2341	2343	2345	2347	2349	2351	2353
2355	2357	2359	2361	2363	2365	2367	2369	2371	2373	2375	2377	2379	2381
2383	2385	2387	2389	2391	2393	2395	2397	2399	2401	2403	2405	2407	2409
2411	2413	2415	2417	2419	2421	2423	2425	2427	2429	2431	2433	2435	2437
2439	2441	2443	2445	2447	2449	2451	2453	2455	2457	2459	2461	2463	2465
2467	2469	2471	2473	2475	2477	2479	2481	2483	2485	2487	2489	2491	2493
2495	2497	2499	2501	2503	2505	2507	2509	2511	2513	2515	2517	2519	2521
2523	2525	2527	2529	2531	2533	2535	2537	2539	2541	2543	2545	2547	2549
2551	2553	2555	2557	2559	2561	2563	2565	2567	2569	2571	2573	2575	2577
2579	2581	2583	2585	2587	2589	2591	2593	2595	2597	2599	2601	2603	2605
2607	2609	2611	2613	2615	2617	2619	2621	2623	2625	2627	2629	2631	2633
2635	2637	2639	2641	2643	2645	2647	2649	2651	2653	2655	2657	2659	2661
2663	2665	2667	2669	2671	2673	2675	2677	2679	2681	2683	2685	2687	2689
2691	2693	2695	2697	2699	2701	2703	2705	2707	2709	2711	2713	2715	2717
2719	2721	2723	2725	2727	2729	2731	2733	2735	2737	2739	2741	2743	2745
2747	2749	2751	2753	2755	2757	2759	2761	2763	2765	2767	2769	2771	2773
2775	2777	2779	2781	2783									

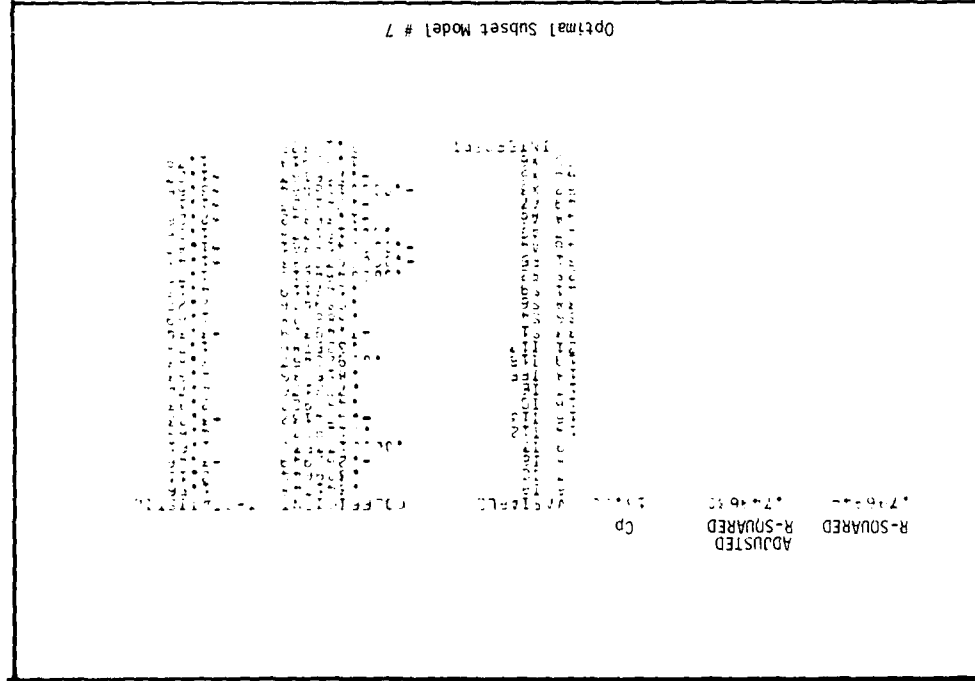
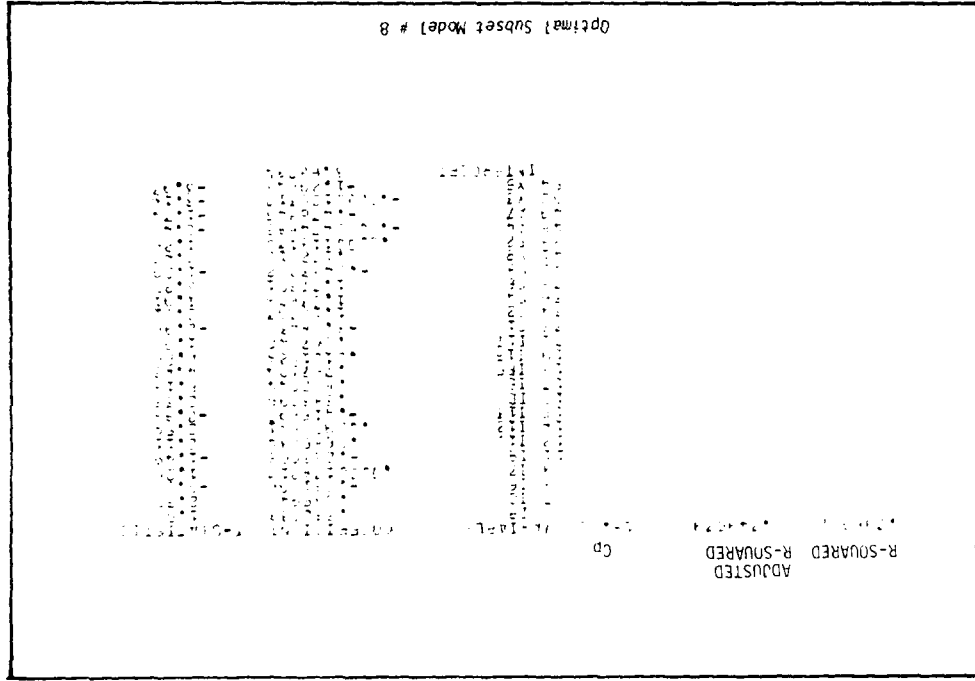


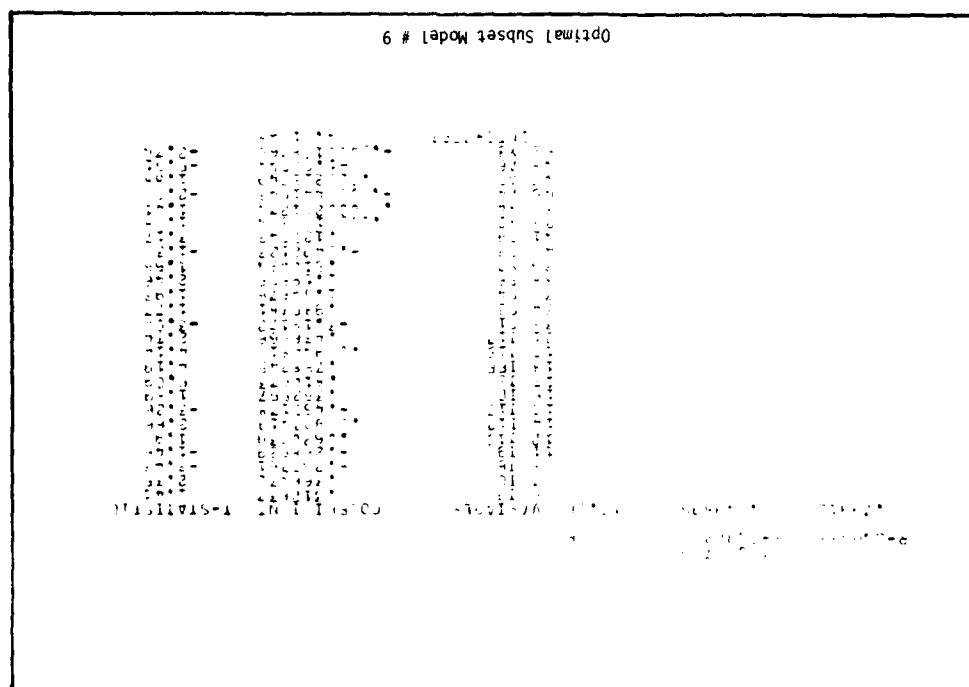
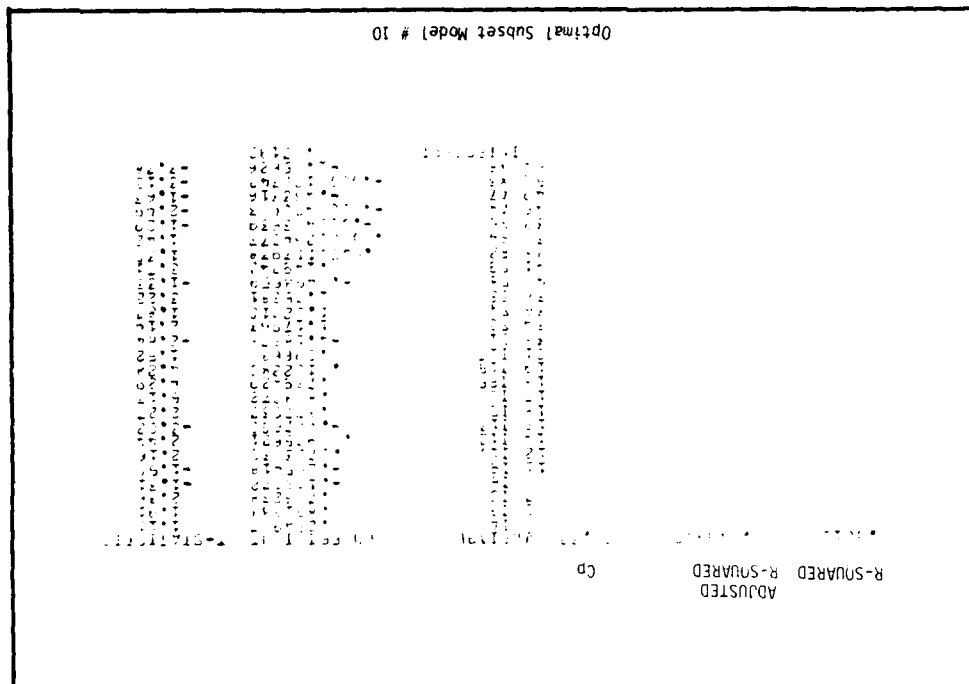
ADDITIONAL INFORMATION  
FROM OUR INFORMATION TO YOU



Optimal Subset Model # 6

Optimal Subset Model # 5





## APPENDIX C

TEN MULTIPLE LINEAR REGRESSION MODELS  
BUILT FROM ESTIMATION SET DATA B

## Variable Key

Factor Label Code/Variable	Factor Description
IA X 1	Dummy Variables Indicating Distributed System Topology
IB X 2	
IC X 3	
16 X 4	Number of Nodes in the Distributed System
17 X 5	
19 X 6	
110 X 7	
112 X 8	
113 X 9	Module to Module Interaction Frequency
JD X10	Dummy Variables Indicating Distribution Policy
IE X11	
JF X12	
115 X13	Percent Nodes Lost
117 X14	
R2 X15	
R4 X16	Initial Assignment Result
R5 X17	Global Memory Capacity
R6 X18	Available Processing Capacity
S1 X19	after Initial Assignment
S2 X20	Available Memory Capacity
S4 X21	after Initial Assignment
S6 X22	Available Communications Capacity
S7 X23	after Initial Assignment
S8 X24	Distributed System Connectivity
	Memory Requirements/Useable Memory Capacity
	Communications Requirements/Useable
	Communications Capacity
	Application System Connectivity
	Dispersion - Initial
	(Number of nodes over which an
	application system is distributed/
	(Number of application system modules)
	Memory Consistency
	(Number of application system modules)/
	(Average number of application system
	modules that will "fit" on a node -
	memory wise)





[illegible]

MULTIPLE LINEAR REGRESSION MODEL # 4									
14	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
15	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
16	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
17	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
18	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
19	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
20	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
21	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
22	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
23	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
24	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
25	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
26	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
27	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
28	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
29	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
30	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
31	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
32	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
33	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
34	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
35	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
36	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
37	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
38	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
39	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
40	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
41	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
42	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
43	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
44	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
45	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
46	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
47	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
48	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
49	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
50	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
51	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
52	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
53	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
54	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
55	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
56	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
57	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
58	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
59	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
60	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
61	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
62	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
63	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
64	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
65	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
66	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
67	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
68	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
69	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
70	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
71	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
72	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
73	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
74	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
75	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
76	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
77	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
78	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
79	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
80	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
81	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
82	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
83	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
84	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
85	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
86	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
87	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
88	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
89	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
90	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
91	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
92	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
93	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
94	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
95	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
96	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
97	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
98	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
99	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
100	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

THIS PAGE IS NOT TO BE REPRODUCED  
FROM COPY FURNISHED TO EDC

Multiple Linear Regression Model # 7

Variable	Parameter	Estimate	Standard Error	t-Statistic	Probability
INTERCEPT		1.0000	.0000		
X1		1.0000	.0000		
X2		1.0000	.0000		
X3		1.0000	.0000		
X4		1.0000	.0000		
X5		1.0000	.0000		
X6		1.0000	.0000		
X7		1.0000	.0000		
X8		1.0000	.0000		
X9		1.0000	.0000		
X10		1.0000	.0000		
X11		1.0000	.0000		
X12		1.0000	.0000		
X13		1.0000	.0000		
X14		1.0000	.0000		
X15		1.0000	.0000		
X16		1.0000	.0000		
X17		1.0000	.0000		
X18		1.0000	.0000		
X19		1.0000	.0000		
X20		1.0000	.0000		
X21		1.0000	.0000		
X22		1.0000	.0000		
X23		1.0000	.0000		
X24		1.0000	.0000		
X25		1.0000	.0000		
X26		1.0000	.0000		
X27		1.0000	.0000		
X28		1.0000	.0000		
X29		1.0000	.0000		
X30		1.0000	.0000		
X31		1.0000	.0000		
X32		1.0000	.0000		
X33		1.0000	.0000		
X34		1.0000	.0000		
X35		1.0000	.0000		
X36		1.0000	.0000		
X37		1.0000	.0000		
X38		1.0000	.0000		
X39		1.0000	.0000		
X40		1.0000	.0000		
X41		1.0000	.0000		
X42		1.0000	.0000		
X43		1.0000	.0000		
X44		1.0000	.0000		
X45		1.0000	.0000		
X46		1.0000	.0000		
X47		1.0000	.0000		
X48		1.0000	.0000		
X49		1.0000	.0000		
X50		1.0000	.0000		
X51		1.0000	.0000		
X52		1.0000	.0000		
X53		1.0000	.0000		
X54		1.0000	.0000		
X55		1.0000	.0000		
X56		1.0000	.0000		
X57		1.0000	.0000		
X58		1.0000	.0000		
X59		1.0000	.0000		
X60		1.0000	.0000		
X61		1.0000	.0000		
X62		1.0000	.0000		
X63		1.0000	.0000		
X64		1.0000	.0000		
X65		1.0000	.0000		
X66		1.0000	.0000		
X67		1.0000	.0000		
X68		1.0000	.0000		
X69		1.0000	.0000		
X70		1.0000	.0000		
X71		1.0000	.0000		
X72		1.0000	.0000		
X73		1.0000	.0000		
X74		1.0000	.0000		
X75		1.0000	.0000		
X76		1.0000	.0000		
X77		1.0000	.0000		
X78		1.0000	.0000		
X79		1.0000	.0000		
X80		1.0000	.0000		
X81		1.0000	.0000		
X82		1.0000	.0000		
X83		1.0000	.0000		
X84		1.0000	.0000		
X85		1.0000	.0000		
X86		1.0000	.0000		
X87		1.0000	.0000		
X88		1.0000	.0000		
X89		1.0000	.0000		
X90		1.0000	.0000		
X91		1.0000	.0000		
X92		1.0000	.0000		
X93		1.0000	.0000		
X94		1.0000	.0000		
X95		1.0000	.0000		
X96		1.0000	.0000		
X97		1.0000	.0000		
X98		1.0000	.0000		
X99		1.0000	.0000		
X100		1.0000	.0000		

Multiple Linear Regression Model # 6

Variable	Parameter	Estimate	Standard Error	t-Statistic	Probability
INTERCEPT		1.0000	.0000		
X1		1.0000	.0000		
X2		1.0000	.0000		
X3		1.0000	.0000		
X4		1.0000	.0000		
X5		1.0000	.0000		
X6		1.0000	.0000		
X7		1.0000	.0000		
X8		1.0000	.0000		
X9		1.0000	.0000		
X10		1.0000	.0000		
X11		1.0000	.0000		
X12		1.0000	.0000		
X13		1.0000	.0000		
X14		1.0000	.0000		
X15		1.0000	.0000		
X16		1.0000	.0000		
X17		1.0000	.0000		
X18		1.0000	.0000		
X19		1.0000	.0000		
X20		1.0000	.0000		
X21		1.0000	.0000		
X22		1.0000	.0000		
X23		1.0000	.0000		
X24		1.0000	.0000		
X25		1.0000	.0000		
X26		1.0000	.0000		
X27		1.0000	.0000		
X28		1.0000	.0000		
X29		1.0000	.0000		
X30		1.0000	.0000		
X31		1.0000	.0000		
X32		1.0000	.0000		
X33		1.0000	.0000		
X34		1.0000	.0000		
X35		1.0000	.0000		
X36		1.0000	.0000		
X37		1.0000	.0000		
X38		1.0000	.0000		
X39		1.0000	.0000		
X40		1.0000	.0000		
X41		1.0000	.0000		
X42		1.0000	.0000		
X43		1.0000	.0000		
X44		1.0000	.0000		
X45		1.0000	.0000		
X46		1.0000	.0000		
X47		1.0000	.0000		
X48		1.0000	.0000		
X49		1.0000	.0000		
X50		1.0000	.0000		
X51		1.0000	.0000		
X52		1.0000	.0000		
X53		1.0000	.0000		
X54		1.0000	.0000		
X55		1.0000	.0000		
X56		1.0000	.0000		
X57		1.0000	.0000		
X58		1.0000	.0000		
X59		1.0000	.0000		
X60		1.0000	.0000		
X61		1.0000	.0000		
X62		1.0000	.0000		
X63		1.0000	.0000		
X64		1.0000	.0000		
X65		1.0000	.0000		
X66		1.0000	.0000		
X67		1.0000	.0000		
X68		1.0000	.0000		
X69		1.0000	.0000		
X70		1.0000	.0000		
X71		1.0000	.0000		
X72		1.0000	.0000		
X73		1.0000	.0000		
X74		1.0000	.0000		
X75		1.0000	.0000		
X76		1.0000	.0000		
X77		1.0000	.0000		
X78		1.0000	.0000		
X79		1.0000	.0000		
X80		1.0000	.0000		
X81		1.0000	.0000		
X82		1.0000	.0000		
X83		1.0000	.0000		
X84		1.0000	.0000		
X85		1.0000	.0000		
X86		1.0000	.0000		
X87		1.0000	.0000		
X88		1.0000	.0000		
X89		1.0000	.0000		
X90		1.0000	.0000		
X91		1.0000	.0000		
X92		1.0000	.0000		
X93		1.0000	.0000		
X94		1.0000	.0000		
X95		1.0000	.0000		
X96		1.0000	.0000		
X97		1.0000	.0000		
X98		1.0000	.0000		
X99		1.0000	.0000		
X100		1.0000	.0000		



## ACKNOWLEDGEMENTS

While the inspiration to begin a dissertation must come from within, the motivation to finish comes from family, friends, advisors and sponsors. To all those who owe deep appreciation for a growing experience.

The author is greatly indebted to her advisor, Professor Richard A. DeMillo, for his interest, insight, guidance and competence and to her doctoral committee: Professor Douglas C. Montgomery, Dr. Charles R. Vick, Professor A. Peter Jensen, and Professor Nancy A. Lynch for the benefit of their constructive comments, suggestions and sense of vision.

A special thanks is given to two friends: Dr. Douglas E. Wrege, whose belief and confidence was always there, and Mr. Clyde G. Roby for his technical talent.

Above all, I am grateful to my husband and our family - to Sam whose love, patience and understanding comforted the tired hours, to my dear son, William, for sensitivity and understanding well beyond his years, and to my sweet daughter, Christine, for her smiles and laughter.

This research was sponsored in part by the U.S. Army Research Office, Research Triangle Park, North Carolina, and the U.S. Army Communications Research and Development Command, Fort Monmouth, New Jersey, under Contract DAA629-79-C-0118. Gratitude is expressed to these organizations for their essential support.

Multiple Linear Regression Model # 10

Variable	Mean	Std. Dev.	Minimum	Maximum
DEPENDANT VA	11.000	1.000	9.000	12.000
INDEPENDANT 1	1.000	1.000	0.000	2.000
INDEPENDANT 2	1.000	1.000	0.000	2.000
INDEPENDANT 3	1.000	1.000	0.000	2.000
INDEPENDANT 4	1.000	1.000	0.000	2.000
INDEPENDANT 5	1.000	1.000	0.000	2.000
INDEPENDANT 6	1.000	1.000	0.000	2.000
INDEPENDANT 7	1.000	1.000	0.000	2.000
INDEPENDANT 8	1.000	1.000	0.000	2.000
INDEPENDANT 9	1.000	1.000	0.000	2.000
INDEPENDANT 10	1.000	1.000	0.000	2.000

ANALYSIS OF VARIANCE

Source	Sum of Squares	df	Mean Square	F	Sig.
Regression	1.000	10	.100	1.000	.999
Residual	1.000	10	.100	1.000	.999
Total	2.000	20	.100	1.000	.999

COEFFICIENTS

Variable	Unstandardized Coefficients	Standardized Coefficients	t	Sig.
(Constant)	11.000		11.000	.000
INDEPENDANT 1	.100	.100	1.000	.999
INDEPENDANT 2	.100	.100	1.000	.999
INDEPENDANT 3	.100	.100	1.000	.999
INDEPENDANT 4	.100	.100	1.000	.999
INDEPENDANT 5	.100	.100	1.000	.999
INDEPENDANT 6	.100	.100	1.000	.999
INDEPENDANT 7	.100	.100	1.000	.999
INDEPENDANT 8	.100	.100	1.000	.999
INDEPENDANT 9	.100	.100	1.000	.999
INDEPENDANT 10	.100	.100	1.000	.999

ALL DATA WERE USED IN THIS ANALYSIS

## VITA

Edith W. Martin was born on June 25, 1945 in Chicago, Illinois. She received her Bachelor of Arts degree in psychology from Lake Forest College in 1967 and Master of Science degree in information and computer science from the Georgia Institute of Technology in 1976. She is married to Professor C. Samuel Martin and has two children, William McNutt Martin, III born December 13, 1971 and Christine Katherine Martin born December 23, 1979. She has been a member of the research and management staff of the Engineering Experiment Station at the Georgia Institute of Technology since 1976. Her professional activities include being a member of the editorial review board of Military Electronics Countermeasures, associate editor of the Journal of Systems and Software, chairman of the Computer Architecture Review Subcommittee of the Electronic Industry Association and member of the Institute of Electrical and Electronic Engineers and Association of Computing Machinery.

## BIBLIOGRAPHY

1. Baron, P., "On Distributed Networks," IEEE Trans. Comm. Tech. COM-12, pp. 1-9, 1964.
2. Beaudry, M. C., "Performance Related Reliability Measures for Computing Systems," Proceedings of the Seventh Annual International Conference on Fault-Tolerant Computing, Los Angeles, Ca., pp. 16-21, June, 1977.
3. Bell System Technical Journal, Vol. 56, No. 7, September, 1977.
4. Borgerson, B. R. and R. F. Freitas, "A Reliability Model for Gracefully Degrading and Standby-sparing Systems," IEEE Trans. on Computers, Vol. C-24, pp. 517-525, May, 1975.
5. Box, G. E. and Hunter, J. S., "The 2 Fractional Factorial Designs Part I," Technometrics, Vol. 3, No. 3, pp. 311-351, August, 1961.
6. Box, G. E. and Hunter, J. S., "The 2 Fractional Factorial Designs Part II," Technometrics, Vol. 3, No. 4, pp. 449-459, November, 1961.
7. DeMillo, R. A., Lipton, R. J., "Software Project Forecasting," School of Information Computer Science, GIT-ICS-80/09, Georgia Institute of Technology, October, 1980.
8. Enslow, P. E., "What is a Distributed Processing System?" Computer, pp. 13-21, January, 1978.
9. Foerster, R. E., "Methodology to Evaluate Strategic Command and Control Systems," Technical Report for HQ USAF, Assistant Chief of Staff Studies and Analysis under Contract No. F44620-74-C-0045, July, 1974.
10. Frank, H. and Frish, I. T., "Analysis and Design of Survivable Networks," IEEE Trans. Comm. Tech., Vol. COM-18, pp. 501-519, October, 1970.
11. Frank, H., "Vulnerability of Communication Networks," IEEE Trans. Comm. Tech., Vol. COM-18, pp. 175-182, September, 1970.
12. Gay, F. A., "Performance Modeling for Distributed Computing Systems," Ph.D. Dissertation, Computer Science Department, Northwestern University, June, 1979.
13. Helmer, O., Reschur, N., "On the Existence of a System of Sciences," Rand Corporation Report No. R-611, December, 1965.

14. Halborn, G., "Measures for distributed processing network survivability," AFIPS Conference Proceedings, National Computer Conference 1980, Vol. 49, pp. 157-163, May, 1980.
15. Jensen, D. E., "The Honeywell Experimental Distributed Processor - An Overview," Computer, pp. 28-37, January, 1978.
16. Losq, J., "A Highly Efficient Redundancy Scheme: Self-Purging Redundancy," IEEE Trans. on Computers, Vol. C-25, pp. 569-578, June, 1976.
17. Losq, J., "Effects of Failure on Gracefully Degradable Systems," Proceedings of the Seventh Annual International Conference on Fault-Tolerant Computing, Los Angeles, Ca., pp. 29-34, June, 1977.
18. Mathur, F. P., Atzlenis, A., "Reliability Analysis and Architecture of a Hybrid-redundant digital system: Generalized triple modular redundancy with self-repair," 1970 SUDC, AFIPS Conference Proceedings, Vol. 36, pp. 375-383, 1970.
19. Merwin, R. E. and Mirhakak, M., "Derivation and Use of a Survivability criterion for DDP systems," AFIPS Conference Proceedings, National Computer Conference 1980, Vol. 49, pp. 139-146, May, 1980.
20. Montgomery, D. C., Peck, E. A., Introduction to Regression Analysis, John Wiley & Sons, Inc., 1981.
21. Montgomery, D. C., "Methods for Factor Screening in Computer Simulation experiments," Technical Report Office of Naval Research, Contract N0014-73-C-0312, March, 1979.
22. Ng, Y. and Avizienis, A., "A Reliability Model for Gracefully Degrading and Repairable Fault-Tolerant Systems," Proceedings of the Seventh Annual International Conference on Fault-Tolerant Computing, Los Angeles, Ca., pp. 22-28, June, 1977.
23. Ng, Y., "Reliability Modeling and Analysis for Fault-Tolerant Computers," Ph.D. Dissertation, Computer Science Department, UCLA, UCLA-ENG-7698, September, 1976.
24. Perlis, A., Sayward, S., Shaw, M., Unpublished notes on software metrics, April, 1980.